

DIPARTIMENTO DI ECONOMIA E FINANZA

METODI E ANALISI STATISTICHE

2020



**UNIVERSITÀ
DEGLI STUDI DI BARI
ALDO MORO**

DIPARTIMENTO DI ECONOMIA E FINANZA

METODI E ANALISI STATISTICHE

2020



UNIVERSITÀ
DEGLI STUDI DI BARI
ALDO MORO

Tutti i diritti di traduzione, riproduzione e adattamento, totale o parziale, con qualsiasi mezzo (comprese le copie fotostatiche e i microfilm) sono riservati

Toma E., d'Ovidio F. (a cura di) (2020). *Metodi e Analisi Statistiche*, Dipartimento di Economia e Finanza, Università degli studi di Bari *Aldo Moro*.

© Copyright 2020 by Università degli Studi di Bari Aldo Moro
www.uniba.it

Prima edizione: dicembre 2020

ISBN 978-88-6629-023-0

Gli articoli qui presentati sono stati oggetto, oltre che di valutazione interna, anche di revisione anonima (in “doppio cieco”).

Editing finale: F.D. d'Ovidio, E. Toma

Sommario

Ernesto Toma Presentazione	pag. 5
Massimo Bilancia, Giovanni Sansaro Inferenza per modelli a topic latenti: una introduzione	« 7
Natalia Leone, Valeria Ancona, Davide Fragnito, Domenico Vitale, Massimo Bilancia Geostatistical analysis of soil reflectance spectra for field-scale digital soil mapping. A case study.	« 95
Claudia Marin, Fabio Manca Modelli generativi per la sentiment analysis	« 119
Davide Fragnito, Natalia Leone, Valeria Ancona, Domenico Vitale, Antonio Lucadamo Comparison of different multivariate calibrations and ensemble methods for estimating selected soil properties with vis-NIR reflectance spectroscopy.	« 135
Maria De Caro, Ilaria Pepe, Paolo Taurisano, Ernesto Toma Studio longitudinale sulle variazioni del (di)stress tra il pre e post-intervento chirurgico per carcinoma mammario in funzione delle strategie di coping adoperate	« 163
Laura Antonucci, Corrado Crocetta, Yana Kostiuk Valutazione dei servizi web offerti dall'Università di Foggia ai tempi del COVID-19	« 189
Laura Antonucci, Corrado Crocetta, Massimo Russo Analisi della soddisfazione rispetto ai servizi erogati da Sanitaservice ASL Foggia s.r.l.	« 205
Nunziata Ribecco, Isabella Stasi, Clelia Punzo, Annalisa Rizzi, Giovanni Tomasicchio Il carcinoma mammario in donne anziane e in donne giovani.	« 231

Salvatore Cariello, Monica Carbonara <i>La povertà nei comuni del Mezzogiorno</i>	« 255
Ezio Ritrovato <i>Crisi economica, prezzi alimentari e delinquenza nella Puglia di fine Ottocento: correlazione e causalità</i>	« 267
Angela Maria D’Uggento, Alessandra Milillo, Barbara Cafarelli <i>È conveniente investire nella finanza sostenibile?</i>	« 281
Agata Maria Madia Carucci, Giovanni Vannella <i>Analisi delle strutture evolutive imprenditoriali a livello regionale e comunale: una applicazione su Puglia e Basilicata</i>	« 293
Crescenzo Gallo, Alessandro Rinaldi <i>Data Envelopment Analysis: application of the efficient frontier on the financial field in the European and the American scenarios</i>	« 323
Mauro Gianfranco Bisceglia <i>La stabilità nelle equazioni differenziali lineari</i>	« 343

Presentazione

Il volume *Metodi e analisi statistiche*, che nel 2020 giunge alla quinta edizione, si conferma un importante momento di confronto per molti ricercatori impegnati sia nell'ambito esclusivamente statistico, che in altri ambiti scientifici che utilizzano il "metodo quantitativo" per analisi matematiche, economiche, storiche, psicologiche, ecc..

Tale metodo diviene in tal senso l'emblema dell'unione di contributi non sempre facilmente classificabili secondo espliciti ed inequivocabili riferimenti a settori scientifico disciplinari. Ci si riferisce, cioè, alle talora troppo rigide distinzioni tra lavori puramente metodologici, ed elaborati prodotti da ambiti diversificati quali la statistica economica, la demografia, la sociologia, la psicologia. Al fine di rendere il volume il più possibile armonico, si è comunque tentato di aggregare le tematiche comuni cercando di rispondere il più possibile a queste esigenze.

Il primo contributo, di Bilancia e Sansaro, rappresenta una interessante e vasta rassegna dei principali risultati esistenti per il modello Latent Dirichlet Allocation (LDA), modello a topic latenti per l'analisi non supervisionata di un corpus di documenti testuali, con i relativi metodi computazionali necessari per la non trattabilità della distribuzione di quello che è a tutti gli effetti un modello Bayesiano gerarchico complesso.

Nel secondo lavoro, di Leone, Bilancia *et al.*, viene presentato un approccio metodologico di tipo multivariato per la mappatura digitale del suolo in un'area di studio del Sud Italia. Gli autori giungono alla conclusione secondo la quale tale approccio sembra essere efficiente e affidabile per mappare la variabilità spaziale del suolo.

Marin e Manca forniscono, invece, una interessante disamina sistematica dei modelli generativi della sentiment analysis, famiglia di metodi statistici ampiamente utilizzati a dati testuali grazie alla generalità e robustezza della metodologia.

Nel quarto contributo, di Vitale *et al.*, si torna a parlare di gestione del suolo. Gli autori, dopo aver sottolineato che tale gestione richiede una corretta valutazione delle proprietà chimiche e fisiche del suolo, spesso ottenuta tramite analisi di laboratorio convenzionali costose e i cui risultati non sono sempre ottenibili in tempi rapidi, propongono un approccio alternativo, più veloce e meno costoso, basato sull'uso di metodologia spettroscopiche e multivariate. Nell'articolo è presente uno studio applicato per stimare alcune proprietà del suolo di base, come sabbia, argilla e carbonio organico, con dati raccolti in un'area agricola del Sud Italia.

Pepe *et al.*, nel quinto lavoro del presente volume, affrontano il ruolo dei fattori stressanti e delle strategie di coping in due fasi cruciali del percorso oncologico, il periodo che precede e segue l'intervento chirurgico da carcinoma mammario della malattia oncologica. La ricerca, realizzata presso il Centro di Senologia multidisciplinare del Policlinico di Bari, con un campione di 135 donne con diagnosi di neoplasia mammaria maligna, ha dimostrato come le variazioni del distress tra il pre e post intervento siano funzione delle strategie di coping adoperate. In particolare, attraverso analisi multivariate, gli autori sono giunti ad evidenziare come le modalità di coping adattive siano predittori

statisticamente significativi della variazione del distress, determinandone l'abbassamento.

Il quinto e sesto contributo hanno come comune denominatore il territorio di riferimento, quello di Foggia. Antonucci, Crocetta e Kostiuk, illustrano i risultati di una indagine sulla efficacia dei servizi web offerti dall'Università di Foggia durante la pandemia da COVID-19. Attraverso l'uso di modelli PLS-PM, gli autori individuano le principali criticità cui dedicare una maggiore attenzione per meglio rispondere alle esigenze degli studenti dell'ateneo dauno. Antonucci, Crocetta e Russo, invece, proseguono gli studi già proposti nella edizione 2019 del presente volume, cercando di individuare i punti di forza e di debolezza per ogni servizio erogato da Sanitaservice ASL FG s.r.l., evidenziando come tale società di servizi abbia sensibilmente migliorato le proprie capacità di problem solving, mantenendo elevati standard qualitativi.

Nell'ottavo contributo si torna ad affrontare la problematica del carcinoma mammario nelle donne. Ribecco et al. confrontano due campioni di donne classificate secondo l'età in "giovani" ed "anziane" al fine di verificare se l'aggressività di sottotipi di carcinoma mammario possa essere influenzata da alcune variabili quali, ad esempio, la dimensione del tumore, il grado istologico dello stesso, il tipo di intervento.

Cariello e Carbonara, dell'ISTAT Puglia, studiano la povertà nei comuni del Mezzogiorno stimando le dimensioni del fenomeno attraverso i dati del Progetto Archivio Integrato Microdati Economici e Demografici (Arch.I.M.E.De). Questi provengono dall'integrazione di basi di dati amministrative e consentono di ottenere microdati relativi alle caratteristiche socio-demografiche e al reddito lordo delle famiglie residenti nei comuni italiani.

Nel decimo contributo, Ritrovato utilizza il "metodo statistico", cui ad inizio di questa presentazione si faceva riferimento, per analizzare le relazioni tra la crisi economica, i prezzi alimentari e la delinquenza nella Puglia di fine Ottocento, fornendo l'ennesima dimostrazione della "contaminazione" statistica negli studi di altri ambiti scientifici.

D'Uggento *et al.* affrontano, invece, l'argomento della "finanza sostenibile" effettuando un confronto tra le banche a prevalente carattere etico e sostenibile (quelle che finanziano in particolare progetti sociali, ambientali e culturali) e l'insieme delle banche europee. L'obiettivo è quello di verificare se le prime mostrino una certa solidità sotto il profilo finanziario-patrimoniale e siano anche dotate di redditività sotto il profilo economico avendo, pertanto, performance simili al complesso delle banche.

Si prosegue con Carucci e Vannella che effettuano una analisi dettagliata delle evoluzioni competitive imprenditoriali, con dettaglio dei dati a livello comunale, utilizzando una applicazione dell'analisi Shift-share su dati derivanti dall'archivio ASIA Unità Locali.

Infine, dopo il contributo in cui Gallo e Rinaldi utilizzano la metodologia DEA (Data Envelopment Analysis) per identificare la frontiera ottimale ed efficiente dei titoli relativi alle aree europee ed americane, Bisceglia fa il punto sul problema della stabilità nelle equazioni differenziali ordinarie di ordine k , realizzando anche una interessante applicazione di un modello economico.

Certi che il volume potrà riscontrare il favore dei lettori, si ringraziano tutti coloro che vi hanno contribuito, sia Autori che referee.

Inferenza per modelli a topic latenti: una introduzione

Massimo Bilancia^{1*}, Giovanni Sansaro²

¹Dipartimento Jonico in “Sistemi Giuridici ed Economici del Mediterraneo: società, ambiente, culture” – Università degli Studi di Bari Aldo Moro, Via Duomo 259 Taranto IT,

²Corso di Laurea Magistrale in Data Science, Dipartimento di Informatica – Università degli Studi di Bari Aldo Moro, Via Orabona 4, 70125 Bari IT

Riassunto: Il modello Latent Dirichlet Allocation (LDA), introdotto in Blei et al. (2003), è il più importante modello a topic latenti per l’analisi non supervisionata di un corpus di documenti testuali. La probabilità di occorrenza di ciascun termine di un documento, pur non dipendendo dalla posizione occupata nel documento, è un miscuglio di multinomiali le cui probabilità di occorrenza dipendono dal topic che genera il particolare token che stiamo considerando. In questo modo, token distinti possono essere generati da topic diversi (non esiste un contenuto tematico globale), e l’identificazione del significato dei topic latenti diventa un potente strumento di riduzione della dimensionalità del corpus di documenti che stiamo studiando.

Da un punto di vista formale il modello LDA è un modello Bayesiano gerarchico complesso. La non trattabilità della relativa distribuzione a posteriori ha reso necessaria l’introduzione di metodi numerici di calcolo ad hoc, come ad esempio quelli che sono basati sull’inferenza variazionale. Altri metodi più tradizionali, come il campionamento di Gibbs, possono essere opportunamente adattati per ridurre la complessità computazionale derivante dalla mole dei dati trasportata da un corpus di documenti testuali. Questi metodi numerici sono tuttavia caratterizzati da una elevata complessità di esposizione. Pertanto, questo paper raccoglie in un contributo unificato i principali risultati esistenti per il modello LDA e i relativi metodi computazionali.

Keywords: apprendimento non-supervisionato; analisi di corpus testuali; latent Dirichlet allocation; modelli Bayesiani gerarchici; inferenza variazionale.

* Autore corrispondente: massimo.bilancia@uniba.it

Entrambi gli autori hanno esaminato e rivisto il manoscritto, contribuendo in egual misura alla stesura, eccetto per la Sezione 11 che è stata scritta interamente da Giovanni Sansaro, e hanno approvato la versione finale accettando di presentare il manoscritto per la pubblicazione.

1. Introduzione

Il problema della **classificazione testuale** è un importante capitolo della teoria dell'apprendimento supervisionato (Sebastiani, 2002). Nella sua forma di base abbiamo a disposizione un corpus di documenti, che deve essere opportunamente pre-processato per essere ridotto ad un insieme di token, ossia di parole che formano un vocabolario dei termini V . Se facciamo l'ipotesi che a ciascun documento del corpus sia associata una etichetta che descrive il contenuto tematico del documento (ad esempio: sport, politica, scienza, religione, etc...), l'obiettivo è quello di costruire un classificatore che sia in grado di predire, con la maggior accuratezza possibile, il contenuto tematico di nuovi documenti che non erano disponibili nel corpus iniziale.

Un tipico esempio di classificatore testuale discriminativo è la **regressione logistica multinomiale** (Ng e Jordan, 2001), nel quale le probabilità a posteriori delle etichette sono apprese sui documenti del corpus utilizzando direttamente il modello, sulla base di una opportuna funzione di un insieme di variabili previsive, calcolate dalle frequenze di occorrenza di ciascun token del vocabolario dei termini. Le probabilità a posteriori possono essere successivamente utilizzate per classificazione dei nuovi documenti che non fanno parte del corpus utilizzato per l'apprendimento.

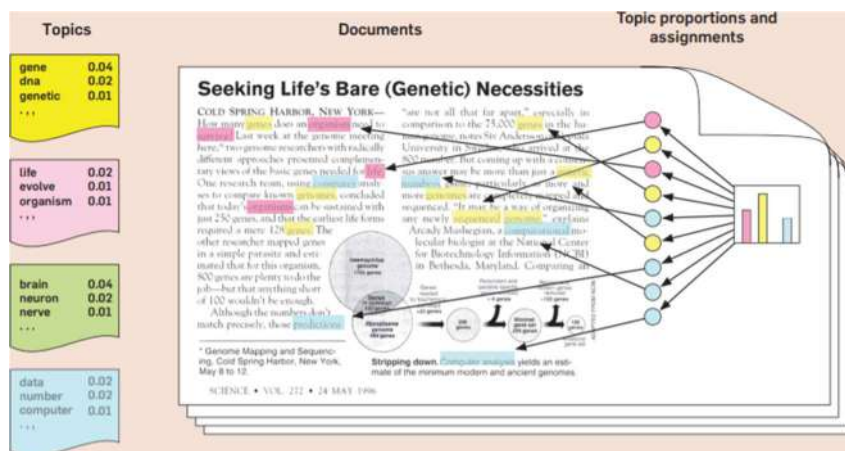
L'idea che la deconvoluzione della semantica di un testo possa essere basata sulla frequenza di occorrenza dei token del vocabolario dei termini è intuitivamente ragionevole, ed è alla base di modelli più elaborati, che si fondano direttamente sull'idea che la probabilità di occorrenza di un token all'interno di un documento non dipende dalla posizione nella quale il token compare. In questa classe, il capostipite è rappresentato dal **modello unigram** (Manning et al., 2008). In esso, la verosimiglianza del modello è basata sulla distribuzione Multinomiale, e dipende dai token solo ed esclusivamente attraverso le relative frequenze di occorrenza nei documenti. Quindi, ciascun testo è trattato come una **bag-of-words** (borsa di parole), nella quale l'ordine di occorrenza dei termini che compongono il test è assolutamente irrilevante.

I problemi di natura supervisionata non esauriscono, però, tutte le possibilità a disposizione. Sempre nell'ipotesi che a ciascun documento possa essere associato in modo univoco un unico contenuto tematico (etichetta), possiamo estendere il modello unigram in modo **non-supervisionato** definendo un miscuglio di modelli unigram, nel quale la probabilità condizionale di occorrenza di ciascun token data l'etichetta è ancora descritta tramite una distribuzione multinomiale definita sul vocabolario dei termini V , e l'insieme delle etichette è descritto da una distribuzione di probabilità discreta (Nigam et al., 2000). In questo caso, il significato delle etichette non

è disponibile a priori, né è noto il numero di etichette, e sarà quindi necessario un ulteriore step di analisi per individuare quale sia il contenuto semantico effettivo che possa essere associato all'etichetta che è stata individuata. In altre parole, una volta appreso il modello, tutto quello che siamo in grado di fare è di partizionare i documenti nelle classi a disposizione, ossia attribuire a ciascun documento una ed una sola etichetta (hard clustering) sulla base delle stime delle probabilità a posteriori delle etichette stesse.

Lo sviluppo più importante in questa classe di modelli è rappresentato dal modello noto come **Latent Dirichlet Allocation (LDA)**, introdotto in Blei et al. (2003) assieme ad un metodo di stima dei parametri ad-hoc particolarmente efficiente dal punto di vista computazionale. Se nell'approccio non-supervisionato possiamo identificare ciascuna etichetta con il contenuto tematico univoco del documento, nel modello LDA è permessa la coesistenza di più contenuti tematici in uno stesso documento, contenuti che prendono il nome di **topic**.

Figura 1. L'intuizione dietro al modello LDA. Ciascun topic è una distribuzione di probabilità sullo stesso vocabolario dei termini $|V|$, e i termini di un documento appartengono a topic distinti, presenti nel documento con una certa frequenza relativa (tratto da Blei, 2012).



La probabilità di occorrenza di ciascun termine di un documento, pur continuando a non dipendere dalla posizione occupata nel documento, è un miscuglio di multinomiali le cui probabilità dipenderanno dal topic che genera il particolare token che stiamo considerando. In questo modo, token distinti possono essere generati da topic diversi (non esiste un contenuto tematico globale), e l'identificazione del significato dei topic latenti diventa un potente strumento di riduzione della dimensionalità del corpus di documenti che stiamo studiando. In altre parole, l'intero corpus può essere

riassunto da un numero finito di topic, ai quali può essere più agevole associare una semantica che ci permette di comprendere il contenuto di un documento nel quale prevalgono alcuni topic invece di altri. L'intuizione che sta dietro il modello LDA può essere meglio colta guardando alla Figura 1, tratta da (Blei, 2012):

Da punto di vista formale il modello LDA è un modello Bayesiano gerarchico complesso (una buona introduzione ai modelli gerarchici è: Carlin e Louis, 2008). La non trattabilità della distribuzione a posteriori ha reso l'introduzione di metodi numerici di calcolo ad hoc, come ad esempio quelli che sono basati sull'**inferenza variazionale** (Blei et al., 2017). Altri metodi più tradizionali, come il campionamento di Gibbs, possono essere opportunamente adattati per ridurre la complessità computazionale derivante dalla mole dei dati trasportata da un corpus testuale.

Questi metodi numerici sono tuttavia caratterizzati da una elevata complessità di trattazione. Abbiamo pertanto deciso di raccogliere tutti i principali risultati relativi allo sviluppo del modello LDA e alla sua trattazione computazionale in un unico paper. Finora, tali risultati, e le relative dimostrazioni, erano disponibili solo parzialmente, e parcellizzati in un certo numero di rapporti tecnici disponibili sul Web (ad esempio: Heinrich, 2008). Queste difficoltà di carattere tecnico potrebbero rallentare la curva di apprendimento del modello per tutti quei ricercatori che volessero utilizzare il modello LDA come strumento teorico da analizzare ed estendere, ovvero come strumento per analizzare corpus testuali reali nella ricerca applicata. Per questi motivi, gli sviluppi formali contenuti in questo lavoro saranno particolarmente dettagliati.

A tale scopo, il lavoro è organizzato come segue. La Sezione 2 richiama, per completezza di trattazione, alcuni strumenti matematici essenziali quali la divergenza di Kullback-Leibler. La Sezione 3 contiene la specificazione del modello LDA e l'espressione dei calcoli dettagliati per arrivare all'espressione della distribuzione a posteriori dei parametri, che non è esplicitabile. La Sezione 4 riguarda il primo algoritmo di stima proposto per aggirare la non trattabilità della distribuzione a posteriori, e cioè l'algoritmo EM variazionale. Nella Sezione 5 descriviamo i dettagli dell'algoritmo EM variazionale quando esso viene utilizzato specificatamente per il modello LDA. Nella Sezione 6 viene trattata una versione leggermente estesa del modello LDA in senso pienamente Bayesiano, descrivendo il relativo algoritmo variazionale. La Sezione 7 è dedicata ad esporre una versione computazionalmente efficiente dell'algoritmo Gibbs Sampling, adattata al modello LDA. Nella Sezione 8 discutiamo la scelta del numero ottimale di topic nell'ambito della scelta Bayesiana tra modelli. La Sezione 9 è dedicata ad una breve rassegna (senza il dettaglio delle parti precedenti) degli sviluppi più recenti. La Sezione 10 illustra le principali librerie

disponibili per la stima dei modelli a topic latenti. Nella Sezione 11 riportiamo un breve caso di studio esemplificativo. Infine, la Sezione 12 raccoglie le fila di questo lavoro, e delinea quelle che sono le aree di ricerca futura più promettenti.

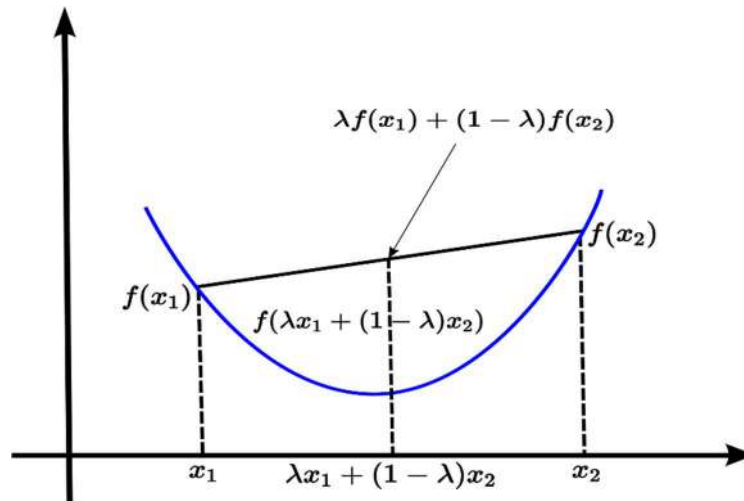
2. Alcuni strumenti matematici

Richiamiamo preliminarmente alcune definizioni e risultati che sono assolutamente necessari per il prosieguo della trattazione. Se $f: I \rightarrow \mathbb{R}$ è una funzione a valori reali definita sull'intervallo I , aperto o chiuso, limitato o illimitato, diremo che f è **convessa** in I se e solo se per ogni $x_1, x_2 \in I$, con $x_1 \neq x_2$, e per ogni $\lambda \in [0,1]$:

$$f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2). \quad (1)$$

La funzione $y = \lambda f(x_1) + (1 - \lambda)f(x_2)$ descrive, al variare di $\lambda \in [0,1]$, un segmento avente per estremi il punto di coordinate $(x_1, f(x_1))$, ottenuto per $\lambda = 0$, e il punto $(x_2, f(x_2))$ corrispondente a $\lambda = 1$. Invece, $f(\lambda x_1 + (1 - \lambda)x_2)$ è l'immagine tramite la funzione f di un punto appartenente al segmento $\lambda x_1 + (1 - \lambda)x_2$ che giace sull'asse delle ascisse. Dunque, in una funzione convessa il grafico giace sempre al di sotto di ogni corda; ecco l'interpretazione geometrica standard per una funzione convessa:

Figura 2. Esempio di funzione convessa.



Infine, diremo che f è **concava** (risp.: strettamente concava) se e solo se $-f$ è convessa (risp.: strettamente convessa).

L'importanza del concetto di convessità per i nostri scopi risiede nel fatto che esso è collegato alla **disuguaglianza di Jensen**: se X è una variabile aleatoria ed f è una funzione convessa sull'intervallo dei valori assunti da X , allora:

$$E[f(X)] \geq f[E(X)], \quad (2)$$

ovvero, nel caso in cui f sia concava:

$$E[f(X)] \leq f[E(X)]. \quad (3)$$

La disuguaglianza di Jensen assume importanza in relazione alla **divergenza di Kullback-Leibler (KL)**, che serve per misurare la 'vicinanza' tra distribuzioni di probabilità. Nel caso discreto, se $p(x)$ e $q(x)$ sono funzioni di probabilità, indichiamo con \mathcal{X} il supporto della distribuzione di probabilità indotta da q , ossia:

$$\mathcal{X} = \text{supp}(q) = \{x | q(x) > 0\}. \quad (4)$$

La divergenza di KL di q da p è il seguente valore atteso:

$$D(p||q) = \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{q(x)} = E_p \left[\log \frac{p(X)}{q(X)} \right]. \quad (5)$$

Se indichiamo con $A = \text{supp}(p) = \{x | p(x) > 0\}$ il supporto della distribuzione di probabilità p , nel calcolo della divergenza di KL dobbiamo tener conto delle seguenti particolarità:

- per convenzione:

$$0 \log \frac{0}{0} = 0,$$

ossia i punti che sono al di fuori del supporto sia di p che di q non contribuiscono al calcolo della divergenza di KL.

- se $A \subset \mathcal{X}$ esistono valori che hanno probabilità non nulla rispetto a q ma nulla rispetto a p , e tali punti non contribuiscono al calcolo della divergenza di KL, poiché per essi per convenzione poniamo:

$$0 \log \frac{0}{q(x)} = 0.$$

- se $A \supset X$ siamo nel caso opposto di quello precedente, ossia esistono valori che hanno probabilità non nulla rispetto a p ma nulla rispetto a q . Per tali valori si pone:

$$p(x) \log \frac{p(x)}{0} = +\infty,$$

e quindi in questo caso la divergenza di KL è automaticamente infinita: $D(p||q) = +\infty$.

Sfruttando la concavità della funzione logaritmo, ed utilizzando la disuguaglianza di Jensen è possibile dimostrare che (si veda, ad esempio, Cover e Thomas, 2005):

- $D(p||q) \geq 0$
- $D(p||q) = 0$ se e solo se $p(x) = q(x)$ per ogni $x \in \mathcal{X}$.

Dunque, il range dei valori assunti da $D(p||q)$ è $[0, +\infty]$, ed affinché la disuguaglianza di KL si mantenga finita è necessario che il supporto di p sia contenuto nel supporto di q , altrimenti la divergenza di KL è automaticamente infinita. Inoltre, tale misura è nulla se e solo $A = \mathcal{X}$ e $p(x) = q(x)$ su ogni punto del supporto comune.

Pertanto, è naturale interpretare $D(p||q)$ come una misura che indica di quanto la distribuzione **approssimata** $p(x)$ differisce dalla distribuzione **approssimante** $q(x)$, quando a p sostituiamo q . Solo nel caso in cui $p \equiv q$ sul supporto comune la divergenza tra le due distribuzioni è nulla; inoltre, se il supporto di q è strettamente contenuto nel supporto di p la divergenza è infinita, in quanto vogliamo escludere tutte quelle situazioni nella quali la distribuzione approssimante q non può generare valori che invece hanno probabilità non nulla sotto la distribuzione p . Infine, la divergenza di KL non è una vera e propria misura di distanza tra distribuzioni di probabilità, poiché si dimostra con semplici controesempi che in generale non è simmetrica negli argomenti, ossia:

$$D(p||q) \neq D(q||p). \quad (6)$$

Nel caso assolutamente continuo, se $f(x)$ e $g(x)$ sono densità di probabilità, la definizione della divergenza di KL si modifica in modo ovvio:

$$D(f||g) = \int_{\mathcal{X}} f(x) \log \frac{f(x)}{g(x)} dx = E_f \left[\log \frac{f(X)}{g(X)} \right], \quad (6)$$

con $\mathcal{X} = \text{supp}(g) = \{x|g(x) > 0\}$. Anche in questo caso si utilizzano le medesime convenzioni e sono valide le stesse proprietà viste nel caso discreto, ossia che la

divergenza di KL è sempre positiva ed è nulla solo quando $f(x) = g(x)$ per ogni $x \in \mathcal{X}$, ed inoltre è finita se e solo se $\text{supp}(f) \subseteq \text{supp}(g)$, mentre $D(f||g) = +\infty$ se e solo se $\text{supp}(f) \not\subseteq \text{supp}(g)$.

Avremo anche bisogno di utilizzare in modo intensivo la **distribuzione di Dirichlet**, indicata come $\text{Dirichlet}_K(\alpha)$, dove l'iperparametro $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_K)$ è un vettore di scalari positivi, $\alpha_i > 0$ per $i = 1, 2, \dots, K$ e $K \geq 2$. Il supporto della distribuzione di Dirichlet è simplex K -dimensionale \mathcal{S}_K , costituito da tutti i vettori $\theta \in \mathbb{R}^K$ che sommano ad 1, e quindi formano una distribuzione di probabilità discreta valida. La densità di probabilità della distribuzione di Dirichlet è la seguente:

$$p(\theta|\alpha) = \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{i=1}^K \theta_i^{\alpha_i-1}, \quad \theta \in \mathcal{S}_K. \quad (7)$$

Insieme alla distribuzione di Dirichlet, avremo bisogno della **coniugazione Dirichlet-Multinomiale**. Supponiamo di considerare il seguente modello Bayesiano gerarchico (che è parte della struttura del modello gerarchico complesso che discuteremo nel prossimo paragrafo):

$$\begin{aligned} \theta &\sim \text{Dirichlet}_K(\alpha) \\ z|\theta &\sim \text{Multinomial}_K(\theta)' \end{aligned} \quad (8)$$

dove $z = (n_1, n_2, \dots, n_K) \in \mathbb{R}^K$ è un vettore di interi che sommano ad una quantità intera prefissata. Dunque, la distribuzione di Dirichlet è la distribuzione a priori per le probabilità Multinomiali che governano la distribuzione di probabilità di z . Ma allora, la distribuzione a posteriori di θ ha espressione:

$$\begin{aligned} p(\theta|z, \alpha) &= \frac{p(\theta, z|\alpha)}{\int p(\theta, z|\alpha)} \propto p(\theta, z|\alpha) = p(z|\theta)p(\theta|\alpha) = \\ &= \prod_{i=1}^K \theta_i^{\alpha_i-1} \prod_{i=1}^K \theta_i^{n_i} = \prod_{i=1}^K \theta_i^{\alpha_i+n_i-1} \propto \\ &\propto \text{Dirichlet}_K(\alpha + z), \end{aligned} \quad (9)$$

dove, per definizione, $\alpha + z \equiv (\alpha_1 + n_1, \alpha_2 + n_2, \dots, \alpha_K + n_K)$. Dunque, la distribuzione a priori e quella a posteriori per il parametro θ hanno la stessa forma funzionale. Nella letteratura Bayesiana si dice che la distribuzione a priori di Dirichlet è coniugata alla verosimiglianza multinomiale (Gelman et al., 2013).

3. Latent Dirichlet allocation

Come abbiamo scritto nell'introduzione, il modello Latent Dirichlet allocation (LDA) è un modello generativo testuale che permette la compresenza, in ciascun documento, di più argomenti tematici (topic): ogni parola di un documento viene assegnata ad un unico topic. Il riferimento fondamentale è Blei et al. (2003). Introduciamo le definizioni fondamentali per trattare il modello:

- **un token** è un item da un vocabolario V , i cui elementi sono indicizzati mediante l'insieme $\{1, 2, \dots, |V|\}$. Ogni token in V verrà rappresentato come un vettore unitario $|V|$ -dimensionale, nel senso che il token il cui indice nel vocabolario v (l'indice associato a ciascun token è arbitrario e tutte le permutazioni degli indici sono equivalenti) verrà rappresentato tramite il vettore \mathbf{w} avente componenti:

$$\begin{aligned} w^v &= 1 \\ w^u &= 0 \quad \text{per } u \neq v. \end{aligned}$$

Dunque, ciascun token è rappresentato in modo univoco mediante un vettore \mathbf{w} della base ortogonale di $\mathbb{R}^{|V|}$, e le componenti di \mathbf{w} sono indicizzate mediante apici;

- **un documento** d è una successione ordinata di N_d token, ossia:

$$\mathbf{w}_d = (w_{d1}, w_{d2}, \dots, w_{dN_d}).$$

In questo caso, poiché abbiamo una successione ordinata, l'indice di w_n punta alla posizione all'interno del documento nel quale il token appare. Ovviamente, se $\pi(\cdot)$ è una qualsiasi permutazione degli interi da 1 a N_d , i documenti $(w_{d1}, w_{d2}, \dots, w_{dN_d})$ e $(w_{d\pi(1)}, w_{d\pi(2)}, \dots, w_{d\pi(N_d)})$ sono distinti, poiché non abbiamo nessun particolare motivo a priori per non considerare distinti due documenti che differiscono solo per l'ordine nel quale i token appaiono. Tuttavia, come vedremo, il meccanismo generativo del modello LDA considera $(w_{d1}, w_{d2}, \dots, w_{dN_d})$ e una sua qualsiasi permutazione come documenti equivalenti, nel senso che tutte le sequenze di token che differiscono solo per l'ordine in cui i token compaiono hanno la stessa probabilità di verificarsi. Diremo anche che ciascun documento è trattato come una **bag-of-words** (una borsa di parole, nella quale le parole sono mescolate alla rinfusa senza tener conto dell'ordine);

- un corpus \mathcal{D} è un insieme di M documenti, ossia:

$$\mathcal{D} = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_M\}.$$

Vediamo ora il meccanismo generativo di un singolo documento (poiché stiamo considerando un singolo documento rimuoviamo il riferimento d). Trattare un modello gerarchico come un modello generativo, significa enfatizzarne la natura di meccanismo aleatorio che è in grado di produrre uno stream di token (ossia un documento): nel finale di questo paragrafo descriveremo con maggior dettaglio le proprietà che sono possedute dai documenti che condividono un meccanismo generativo probabilistico come quello che abbiamo appena descritto. Naturalmente, nell'uso inferenziale i dati (ossia i documenti) sono già disponibili, e l'approccio viene rovesciato, poiché l'obiettivo è quello di effettuare l'inferenza sui parametri del modello.

Di fatto, il modello LDA è un modello Bayesiano gerarchico avente la seguente struttura:

1. generiamo un vettore K -dimensionale:

$$\theta \sim \text{Dirichlet}_K(\alpha).$$

2. per ciascuno degli N termini del documento ($n = 1, 2, \dots, N$) ed indipendentemente l'uno dall'altro:
 - a. generiamo un topic $z_n | \theta \sim \text{Multinomial}_K(\theta)$;
 - b. generiamo un token dalla distribuzione Multinomiale $|V|$ -dimensionale $p(w_n | z_n, \beta)$, cui parametri sono condizionati al topic z_n che è stato generato.

Il numero di topic K è prefissato e noto in anticipo (la scelta del valore ottimale di K verrà trattata più avanti). La distribuzione Multinomiale utilizzata per generare z_n è una distribuzione Multinomiale su una singola prova. Pertanto, ogni topic viene rappresentato come un vettore unitario della base ortogonale di \mathbb{R}^K .

Poiché la distribuzione Multinomiale $p(w_n | z_n, \beta)$ dalla quale sono generati i token dipende dall'indicatore z_n del topic, di fatto avremo K distinte distribuzioni di probabilità che governano la generazione dei token. Ciò può essere meglio compreso se esplicitiamo il fatto che il parametro β deve essere visto come una matrice $\beta \in \mathbb{R}^{K \times |V|}$, il cui elemento generico è il seguente:

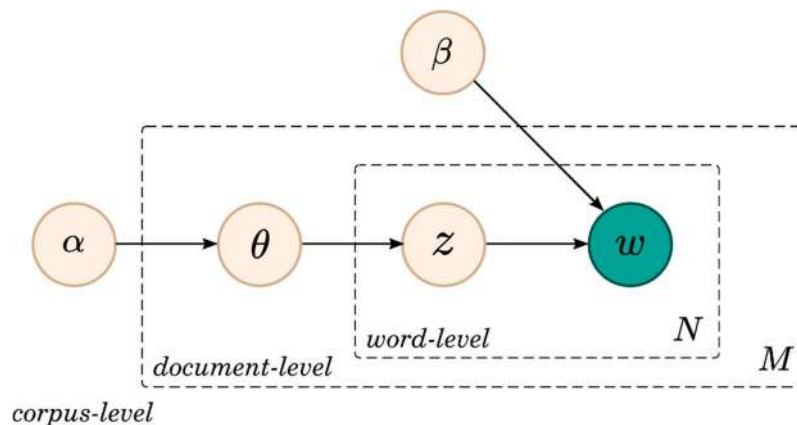
$$\beta_{ij} = Pr(w_n^j = 1 | z_n^i = 1). \quad (10)$$

Utilizzando la stessa notazione basata su apici vista per i token, l'indice di riga viene determinato dall'indicatore del topic, e quindi ogni riga di tale matrice (in tutto abbiamo K righe) è una distribuzione di probabilità sul vocabolario V , che governa la distribuzione di probabilità Multinomiale $|V|$ -dimensionale che presiede alla

generazione dei token del documento. Inoltre, token distinti possono essere generati da topic differenti, e quindi più topic possono coesistere nello stesso documento; anzi, ciascun documento è un miscuglio di topic, e le frequenze di occorrenza di ciascuno dei K sono contenute nel vettore $\theta \in \mathbb{R}^K$ poiché, come sappiamo, le realizzazioni della distribuzione di Dirichlet consistono in distribuzioni di probabilità discrete.

Di fatto, sebbene nel modello un topic z_n è un indicatore K -dimensionale, esso può essere identificato con la distribuzione di probabilità sul vocabolario V ad esso associata. Con K topic avremo K distinte distribuzioni di probabilità che governano la generazione dei token; in generale, uno stesso token avrà probabilità di occorrenza distinte sotto topic differenti (poiché è naturale che certi termini compaiano più frequentemente quando il documento tratta di certe tematiche, e meno frequentemente quando le tematiche sono altre). Infine, notiamo che le variabili indicatrici z_n dei topic non sono direttamente osservabili, motivo per il quale esse sono anche considerate come variabili latenti. La rappresentazione del modello LDA in forma di grafo orientato è contenuta nella Figura 3 sottostante:

Figura 3. Il modello LDA in forma grafica.



Gli iperparametri α (che governano la frequenza di occorrenza dei topic in ciascun documento) e β (le distribuzioni di probabilità sui topic) sono **corpus-wise** (o globali). Invece, come è ovvio, abbiamo un θ distinto per ciascuno degli M documenti del corpus, ossia θ è **document-wise** (o locale). Dunque, se utilizziamo l'indicatore del documento il parametro θ andrebbe riscritto come θ_d , mentre α e β non dipendono in ogni caso dall'indice del documento d .

Osserviamo ora, utilizzando l'espressione della distribuzione Multinomiale, che la verosimiglianza su un singolo termine è:

$$p(w_n|z_n, \beta) = \prod_{j=1}^V \beta_{ij}^{w_n^j}. \quad (11)$$

Data l'ipotesi di indipendenza condizionale che abbiamo fatto, la probabilità di generazione dell'intero documento è pari a:

$$p(\mathbf{w}|\mathbf{z}, \beta) = \prod_{n=1}^N p(w_n|z_n, \beta) = \prod_{n=1}^N \prod_{j=1}^{|V|} \beta_{ij}^{w_n^j}, \quad (12)$$

con $\mathbf{z} = (z_1, z_2, \dots, z_N)$. Naturalmente, l'indicatore i del topic varia termine per termine, ed una notazione più corretta dovrebbe includere esplicitamente il riferimento al termine, ma per non appesantire troppo la notazione eviteremo di includere tale riferimento, dandolo per scontato.

Per quanto riguarda l'indicatore dei topic, abbiamo ovviamente:

$$p(z_n|\theta) = \prod_{i=1}^K \theta_i^{z_n^i}. \quad (13)$$

A questo punto possiamo scrivere agevolmente la probabilità di occorrenza congiunta per il termine in posizione n e il relativo topic:

$$p(w_n, z_n|\beta, \theta) = p(w_n|z_n, \beta)p(z_n|\theta) = \prod_{j=1}^{|V|} \beta_{ij}^{w_n^j} \prod_{s=1}^K \theta_s^{z_n^s}. \quad (14)$$

Per ottenere la distribuzione marginale di w_n , marginalizzando rispetto a z_n otteniamo:

$$p(w_n|\beta, \theta) = \sum_{z_n} p(w_n, z_n|\beta, \theta) = \sum_{z_n} p(w_n|z_n, \beta)p(z_n|\theta), \quad (15)$$

ossia tale probabilità marginale è un miscuglio delle K possibili probabilità multinomiali di occorrenza, ponderate per le probabilità di occorrenza di ciascun topic.

Esplicitare tale probabilità marginale è abbastanza agevole, poiché se $z_n^i = 1$ l'unico termine non uguale ad 1 in:

$$\prod_{s=1}^K \theta_s^{z_n^s},$$

è appunto quello corrispondente ad $s = i$, e quindi da:

$$p(w_n | \beta, \theta) = \sum_{i=1}^K \prod_{j=1}^{|V|} \beta_{ij}^{w_n^j} \prod_{s=1}^K \theta_s^{z_n^s}, \quad (16)$$

segue subito che:

$$p(w_n | \beta, \theta) = \sum_{i=1}^K \prod_{j=1}^{|V|} (\theta_i \beta_{ij})^{w_n^j}. \quad (17)$$

Si noti che è lecito scrivere $\theta_i \times \beta_{ij}^{w_n^j} = (\theta_i \beta_{ij})^{w_n^j}$, poiché w_n^j può essere solo 0 oppure 1. A questo punto, è immediato che la probabilità marginale di occorrenza dell'intero documento è pari a:

$$p(\mathbf{w} | \beta, \theta) = \prod_{n=1}^N p(w_n | \beta, \theta) = \prod_{n=1}^N \sum_{i=1}^K \prod_{j=1}^{|V|} (\theta_i \beta_{ij})^{w_n^j}. \quad (18)$$

Per l'inferenza a posteriori sui parametri di allocazione e le variabili latenti del modello, utilizzando il teorema di Bayes otteniamo:

$$p(\theta, \mathbf{z} | \mathbf{w}, \alpha, \beta) = \frac{p(\theta, \mathbf{z}, \mathbf{w} | \alpha, \beta)}{p(\mathbf{w} | \alpha, \beta)}. \quad (19)$$

Come è standard nell'inferenza Bayesiana, la distribuzione a posteriori è condizionata ai dati e agli iperparametri α e β . Per il numeratore non abbiamo nessun problema, in quanto:

$$\begin{aligned}
p(\theta, \mathbf{z}, \mathbf{w}|\alpha, \beta) &= p(\theta|\alpha)p(\mathbf{z}, \mathbf{w}|\beta, \theta) = \\
&= p(\theta|\alpha) \prod_{n=1}^N p(w_n, z_n|\beta, \theta) = \\
&= p(\theta|\alpha) \prod_{n=1}^N p(w_n|z_n, \beta)p(z_n|\theta),
\end{aligned} \tag{20}$$

e in questa espressione $p(\theta|\alpha)$ può essere immediatamente scritta in modo esplicito poiché $\theta \sim \text{Dirichlet}_K(\alpha)$, mentre l'altro fattore è stato già esplicitato. Ai fini del calcolo della esplicitazione della distribuzione a posteriori, la parte problematica è il denominatore (la **verosimiglianza marginale** o **evidenza** del modello). Infatti:

$$\begin{aligned}
p(\mathbf{w}|\alpha, \beta) &= \int p(\mathbf{w}, \theta|\alpha, \beta) d\theta = \\
&= \int p(\mathbf{w}|\theta, \alpha, \beta)p(\theta|\alpha, \beta) d\theta = \\
&= \int p(\mathbf{w}|\beta, \theta)p(\theta|\alpha) d\theta = \\
&= \int \left[\prod_{n=1}^N p(w_n|\beta, \theta) \right] p(\theta|\alpha) d\theta = \\
&= \int \left[\prod_{n=1}^N \sum_{i=1}^K \prod_{j=1}^{|V|} (\theta_i \beta_{ij})^{w_n^j} \right] p(\theta|\alpha) d\theta = \\
&= \frac{\Gamma(\sum_i \alpha_i)}{\prod_i \Gamma(\alpha_i)} \int \left[\prod_{n=1}^N \sum_{i=1}^K \prod_{j=1}^{|V|} (\theta_i \beta_{ij})^{w_n^j} \right] \left(\prod_{i=1}^K \theta_i^{\alpha_i-1} \right) d\theta,
\end{aligned} \tag{21}$$

che è chiaramente intrattabile. Dunque, la distribuzione a posteriori non è esplicitabile, ed abbiamo bisogno di ricorrere di metodi numerici ad-hoc per valutarla.

Prima di passare alla trattazione di tali metodi numerici, è interessante osservare una proprietà del modello LDA, che era stata già introdotta all'inizio di questa sezione. Osserviamo che dalla precedente catena di uguaglianze è possibile estrapolare:

$$p(\mathbf{w}|\alpha, \beta) = \int \left[\prod_{n=1}^N p(w_n|\beta, \theta) \right] p(\theta|\alpha) d\theta. \tag{22}$$

Ricordiamo ora che una successione di variabili aleatorie (u_1, u_2, \dots, u_N) si dice **scambiabile** se e solo se (Gelman et al., 2013):

$$p(u_1, u_2, \dots, u_N) = p(u_{\pi(1)}, u_{\pi(2)}, \dots, u_{\pi(N)}), \quad (23)$$

per qualsiasi permutazione $\pi(\cdot)$ degli interi $\{1, 2, \dots, N\}$. In questo caso abbiamo, pertanto, una simmetria dal punto di vista probabilistico, nel senso che qualsiasi permutazione dell'ordine delle variabili contenute nel modello genera sequenze che hanno sempre la stessa probabilità di verificarsi. Allo stesso modo diremo che una successione infinita è **infinitamente scambiabile** se e solo se ogni sotto-successione finita è scambiabile nel senso che abbiamo specificato (Kallenberg, 2005). Mediante controesempi si può facilmente dimostrare che una successione infinita di variabili aleatorie indipendenti e identicamente distribuite è anche scambiabile, ma il viceversa non è necessariamente vero. La scambiabilità è pertanto un concetto meno restrittivo, e richiede solo che le variabili aleatorie coinvolte siano simmetriche dal punto di vista del loro comportamento probabilistico.

La nozione di infinita scambiabilità si adatta perfettamente alle nostre necessità: infatti un meccanismo generativo infinitamente scambiabile è ogni meccanismo in grado di produrre uno stream potenzialmente infinito di token, tale che ogni sotto-successione finita (ossia un testo formato da N token) sia scambiabile. Ma in questo contesto, dire che ogni successione finita di ogni N token è scambiabile vuole dire che l'ordine in cui si presentano i token è irrilevante, e quindi il documento può essere visto come una bag-of-words. Tutte le possibili sequenze ottenute per permutazione hanno la medesima probabilità di occorrenza, e dunque l'inferenza a posteriori non è influenzata da quale particolare sequenza sia stata osservata.

Ma come possiamo dimostrare che il meccanismo generativo del modello LDA è infinitamente scambiabile? A tal fine ci viene in aiuto il **teorema di De Finetti** (Gelman et al., 2013), grazie al quale una successione è infinitamente scambiabile se e solo se per ogni sotto-successione finita vale la seguente rappresentazione:

$$p(u_1, u_2, \dots, u_N) = \int \prod_{n=1}^N p(u_n | \phi) \pi(\phi) d\phi, \quad (24)$$

per una opportuna distribuzione di probabilità a priori $\pi(\phi)$. In altre parole, ogni successione scambiabile può essere rappresentata come un miscuglio di variabili aleatorie condizionalmente i.i.d. dato un parametro ϕ , prendendo come funzione peso una opportuna distribuzione di probabilità a priori su ϕ . Ma se si guarda

l'espressione (22) che abbiamo scritto per $p(\mathbf{w}|\alpha, \beta)$, ne segue che la rappresentazione (24) sussiste per il modello LDA, in particolare prendendo $p(w_n|\beta, \theta)$ come variabili aleatorie condizionalmente i.i.d., e $p(\theta|\alpha)$ come funzione peso del miscuglio. Dunque, il modello LDA tratta i testi come delle bag-of-words.

4. Inferenza variazionale

Per il problema del calcolo della distribuzione a posteriori (19) un primo approccio è soltanto parzialmente Bayesiano, nel senso che gli iperparametri α e β sono considerati incogniti, ma ad essi non è assegnata una ulteriore distribuzione di probabilità a priori completamente specificata. Invece, le stime vengono ottenute massimizzando la log-verosimiglianza marginale rispetto ad α e β :

$$\ell(\alpha, \beta) = \sum_{d=1}^M \log p(\mathbf{w}_d|\alpha, \beta), \quad (25)$$

dove $\mathcal{D} = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_M\}$. Le stime ottenute in questo modo prendono il nome di **stime Empirical Bayes**, in quanto il modello non è specificato in senso pienamente Bayesiano, ma gli iperparametri α e β sono stimati direttamente dai dati invece di essere settati ad un valore predefinito (ovvero di essere ulteriormente trattati attraverso opportune distribuzioni a priori di probabilità che aggiungerebbero un ulteriore livello al modello gerarchico).

Nel caso del modello LDA questo approccio non è possibile in maniera diretta, poiché sappiamo che la verosimiglianza marginale non è esplicitabile. L'idea alla base dei metodi variazionali è allora quella di surrogare la distribuzione a posteriori $p(\theta, \mathbf{z}|\mathbf{w}, \alpha, \beta)$ (non esplicitabile, in quanto non è esplicitabile la verosimiglianza marginale) con una **distribuzione variazionale**:

$$q(\theta, \mathbf{z}|\nu), \quad (26)$$

che dipende dagli **iperparametri variazionali** ν . L'approssimazione di $p(\theta, \mathbf{z}|\mathbf{w}, \alpha, \beta)$ con $q(\theta, \mathbf{z}|\nu)$ viene resa il più accurata possibile minimizzando la divergenza di KL di $p(\theta, \mathbf{z}|\mathbf{w}, \alpha, \beta)$ da $q(\theta, \mathbf{z}|\nu)$:

$$\text{KL}(q(\theta, \mathbf{z}|\nu)||p(\theta, \mathbf{z}|\mathbf{w}, \alpha, \beta)). \quad (27)$$

Si noti che la distribuzione ‘approssimante’ $q(\theta, \mathbf{z}|\nu)$ compare a destra nell’espressione della divergenza di KL, motivo per il quale è comune dire che nell’inferenza variazionale viene utilizzata la versione **reverse** della divergenza di KL, mentre la versione standard $\text{KL}(p(\theta, \mathbf{z}|\mathbf{w}, \alpha, \beta)||q(\theta, \mathbf{z}|\nu))$ prende il nome di **forward-KL**.

Per spiegare ulteriormente la natura del metodo, osserviamo che la log-verosimiglianza marginale di una singola osservazione (ossia di un singolo documento, sopprimendo l’indice d) può essere riscritta come:

$$\begin{aligned}
\log p(\mathbf{w}|\alpha, \beta) &= \log \int \sum_{\mathbf{z}} p(\theta, \mathbf{z}, \mathbf{w}|\alpha, \beta) d\theta = \\
&= \log \int \sum_{\mathbf{z}} \frac{p(\theta, \mathbf{z}, \mathbf{w}|\alpha, \beta)}{q(\theta, \mathbf{z}|\nu)} q(\theta, \mathbf{z}|\nu) d\theta \stackrel{\text{Jensen}}{\geq} \\
&\geq \int \log \sum_{\mathbf{z}} \frac{p(\theta, \mathbf{z}, \mathbf{w}|\alpha, \beta)}{q(\theta, \mathbf{z}|\nu)} q(\theta, \mathbf{z}|\nu) d\theta \stackrel{\text{Jensen}}{\geq} \\
&\geq \int \sum_{\mathbf{z}} q(\theta, \mathbf{z}|\nu) \log \left[\frac{p(\theta, \mathbf{z}, \mathbf{w}|\alpha, \beta)}{q(\theta, \mathbf{z}|\nu)} \right] d\theta = \\
&\geq \int \sum_{\mathbf{z}} [\log p(\theta, \mathbf{z}, \mathbf{w}|\alpha, \beta)] q(\theta, \mathbf{z}|\nu) d\theta - \\
&\quad - \int \sum_{\mathbf{z}} [\log q(\theta, \mathbf{z}|\nu)] q(\theta, \mathbf{z}|\nu) d\theta = \\
&= E_q[\log p(\theta, \mathbf{z}, \mathbf{w}|\alpha, \beta)] - \\
&\quad - E_q[\log q(\theta, \mathbf{z}|\nu)] \stackrel{\text{def}}{=} \\
&\stackrel{\text{def}}{=} \mathcal{L}(\nu|\alpha, \beta),
\end{aligned} \tag{28}$$

dove abbiamo utilizzato in modo decisivo la disuguaglianza di Jensen. La quantità $\mathcal{L}(\nu|\alpha, \beta)$, nella quale abbiamo evidenziato la dipendenza tanto dagli iperparametri variazionali quanto dagli iperparametri originali α e β , prende il nome di **energia libera**. Poiché per una arbitraria distribuzione variazionale q abbiamo dimostrato che:

$$\log p(\mathbf{w}|\alpha, \beta) \geq \mathcal{L}(\nu|\alpha, \beta), \tag{29}$$

l’energia libera è un minorante (lower bound) della log-verosimiglianza marginale. Si noti che l’energia libera è funzione tanto dei parametri variazionali quanto degli

iperparametri (oltre che dei dati \mathbf{w}). Spesso, l'energia libera viene chiamata in modo alternativo **Evidence Lower Bound (ELBO)**:

$$\text{ELBO}(q) \equiv \mathcal{L}(v|\alpha, \beta), \quad (30)$$

mettendo in evidenza in modo esplicito la dipendenza dalla distribuzione variazionale. Per completare la spiegazione del metodo osserviamo che, altrettanto agevolmente, è possibile dimostrare che:

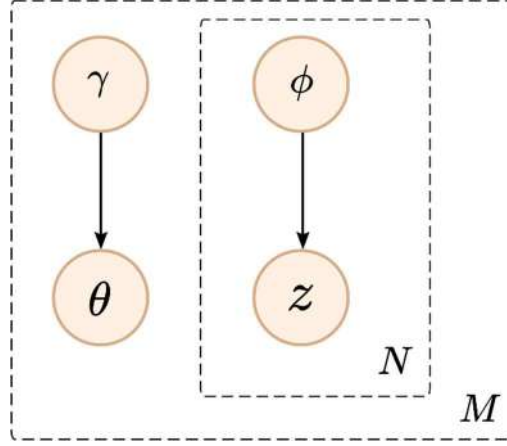
$$\begin{aligned} \log p(\mathbf{w}|\alpha, \beta) &= \mathcal{L}(v|\alpha, \beta) + \\ &+ \text{KL}(q(\theta, \mathbf{z}|v)||p(\theta, \mathbf{z}|\mathbf{w}, \alpha, \beta)). \end{aligned} \quad (31)$$

Dunque, poiché la divergenza di KL è sempre positiva, minimizzare $\text{KL}(q(\theta, \mathbf{z}|v)||p(\theta, \mathbf{z}|\mathbf{w}, \alpha, \beta))$ rispetto ai parametri variazionali significa massimizzare l'energia libera $\mathcal{L}(v|\alpha, \beta)$ rispetto ai parametri variazionali. Se al termine di questa procedura l'approssimazione della distribuzione a posteriori mediante la distribuzione variazionale diventa perfetta, il termine contenente la divergenza di KL diventa nullo. Altrimenti, l'energia libera massimizzata rispetto ai parametri variazionali diventa un surrogato della log-verosimiglianza (in quanto funzione di α e β), e può essere utilizzata per ricavare delle stime EB approssimate di α e β massimizzando l'energia libera rispetto a tali variabili (i parametri variazionali sono stati massimizzati e, di fatto, non sono più quantità variabili).

Per il modello LDA la distribuzione variazionale approssimante è, di solito, di tipo mean-field, ossia:

$$q(\theta, \mathbf{z}|v) = q(\theta, \mathbf{z}|\gamma, \phi) = q(\theta|\gamma) \prod_{n=1}^N q(z_n|\phi_n), \quad (32)$$

Pertanto, in questa approssimazione variazionale le componenti sono tutte a priori indipendenti, e θ e \mathbf{z} sono stati disaccoppiati nel modo seguente:

Figura 4. Approssimazione variazionale della distribuzione a posteriori del modello LDA

Nello specifico, la scelta di ciascuna componente della distribuzione variazionale è la seguente:

$$\begin{aligned} \theta|\gamma &\sim \text{Dirichlet}_K(\gamma), \\ z_n|\phi_n &\stackrel{\text{ind.}}{\sim} \text{Multinomial}_K(\phi_n). \end{aligned} \quad (33)$$

Vediamo, nello specifico, l'espressione assunta dall'energia libera con questa scelta della distribuzione variazionale. Sviluppiamo il primo valore atteso che compone l'energia libera:

$$\begin{aligned} E_q [\log p(\theta, \mathbf{z}, \mathbf{w}|\alpha, \beta)] &= \\ &= E_q [\log \{p(\theta|\alpha)p(\mathbf{z}|\theta)p(\mathbf{w}|\mathbf{z}, \beta)\}] = \\ &= E_q \log[p(\theta|\alpha)] + E_q \log[p(\mathbf{z}|\theta)] + \\ &+ E_q \log[p(\mathbf{w}|\mathbf{z}, \beta)]. \end{aligned} \quad (34)$$

Allo stesso modo, tenendo presente le fattorizzazioni che abbiamo imposto per la distribuzione variazionale, il secondo valore atteso si trasforma in:

$$\begin{aligned} E_q [\log q(\theta, \mathbf{z}|\nu)] &= E_q [q(\theta|\gamma)q(\mathbf{z}|\phi)] = \\ &= E_q \log[q(\theta|\gamma)] + \\ &= E_q \log[q(\mathbf{z}|\phi)], \end{aligned} \quad (35)$$

con $\phi = (\phi_1, \phi_2, \dots, \phi_N)$. Mettendo tutto insieme otteniamo che:

$$\begin{aligned}
\mathcal{L}(\gamma, \boldsymbol{\phi} | \boldsymbol{\alpha}, \boldsymbol{\beta}) &= E_q \log[p(\boldsymbol{\theta} | \boldsymbol{\alpha})] + \\
&+ E_q \log[p(\mathbf{z} | \boldsymbol{\theta})] + \\
&+ E_q \log[p(\mathbf{w} | \mathbf{z}, \boldsymbol{\beta})] - \\
&- E_q \log[q(\boldsymbol{\theta} | \gamma)] - \\
&- E_q \log[q(\mathbf{z} | \boldsymbol{\phi})],
\end{aligned} \tag{36}$$

e possiamo pertanto esplicitare gli addendi che compaiono nella nuova espressione dell'energia libera che abbiamo appena scritto:

1. $E_q \log[p(\boldsymbol{\theta} | \boldsymbol{\alpha})]$.

Ricordiamo innanzitutto che:

$$p(\boldsymbol{\theta} | \boldsymbol{\alpha}) = \frac{\Gamma(\sum_{j=1}^K \alpha_j)}{\prod_{i=1}^K \Gamma(\alpha_i)} \theta_1^{\alpha_1-1} \theta_2^{\alpha_2-1} \dots \theta_K^{\alpha_K-1}, \tag{37}$$

con $\alpha_i > 0$ per ogni i , e quindi:

$$\begin{aligned}
\log p(\boldsymbol{\theta} | \boldsymbol{\alpha}) &= \log \Gamma\left(\sum_{j=1}^K \alpha_j\right) - \sum_{i=1}^K \log \Gamma(\alpha_i) + \\
&+ \sum_{i=1}^K (\alpha_i - 1) \log \theta_i,
\end{aligned} \tag{38}$$

dalla quale segue immediatamente che:

$$\begin{aligned}
E_q \left[\log \Gamma\left(\sum_{j=1}^K \alpha_j\right) \right] &= \log \Gamma\left(\sum_{j=1}^K \alpha_j\right) \\
E_q \left[\sum_{i=1}^K \log \Gamma(\alpha_i) \right] &= \sum_{i=1}^K \log \Gamma(\alpha_i) \\
E_q \left[\sum_{i=1}^K (\alpha_i - 1) \log \theta_i \right] &= \sum_{j=1}^K (\alpha_j - 1) E_q[\log \theta_j].
\end{aligned} \tag{39}$$

In quest'ultima espressione resta da valutare solo il valore atteso $E_q[\log \theta_i]$ rispetto alla distribuzione variazionale q . Naturalmente, quando andiamo ad effettuare i calcoli la densità multinomiale prodotto presente in q scompare, poiché il relativo integrale si fattorizza rispetto a quello in cui compare θ e diventa pari ad 1. Pertanto:

$$E_q[\log \theta_i] = \int \log \theta_i q(\theta|\gamma) d\theta. \quad (40)$$

Per esplicitare tale valore atteso abbiamo bisogno del concetto di **famiglia esponenziale naturale**, che consiste in una famiglia parametrica di distribuzioni di probabilità dipendenti dal parametro naturale $\eta \in \mathbb{R}^K$, avente la seguente espressione:

$$p(\theta|\eta) = h(\theta) \exp\{\eta^\top T(\theta) - A(\eta)\}, \quad (41)$$

dove $T(\eta) = (T_1(\eta), T_2(\eta), \dots, T_K(\eta))$ è il vettore delle **statistiche sufficienti minimali**. La funzione $A(\eta)$ si chiama funzione di partizione: per le famiglie esponenziali naturali abbiamo un notevole risultato che ci permette di ottenere il valore atteso delle statistiche sufficienti minimali utilizzando la funzione di partizione, e cioè:

$$E_p[T_i(\theta)] = \frac{\partial A(\eta)}{\partial \eta_i}, \quad i = 1, 2, \dots, K. \quad (42)$$

La distribuzione di Dirichlet $q(\theta|\gamma)$ è una famiglia esponenziale naturale, come è immediato verificare, poiché possiamo riscriverla come:

$$\begin{aligned} q(\theta|\gamma) &= \exp[\log q(\theta|\gamma)] = \exp\left\{\sum_{i=1}^K (\gamma_i - 1) \log \theta_i - \right. \\ &\quad \left. - \left[\sum_{i=1}^K \log \Gamma(\gamma_i) - \log \Gamma\left(\sum_{j=1}^K \gamma_j\right) \right] \right\}, \end{aligned} \quad (43)$$

con (per $i = 1, 2, \dots, K$):

$$\begin{aligned} \log \theta_i &\rightarrow \text{statistiche sufficienti minimali,} \\ \eta_i &= \gamma_i - 1 \rightarrow \text{parametro naturale,} \\ A(\eta) &= \sum_{i=1}^K \log \Gamma(\eta_i + 1) - \log \Gamma\left[\sum_{j=1}^K (\eta_j + 1)\right]. \end{aligned} \quad (44)$$

Pertanto:

$$\begin{aligned}
E_q[\log \theta_i] &= \frac{\partial A(\eta)}{\partial \eta_i} = \\
&= \frac{\partial A(\gamma_i)}{\partial \gamma_i} \frac{\partial \gamma_i}{\partial \eta_i} = \\
&= \frac{\partial A(\gamma_i)}{\partial \gamma_i} \frac{\partial(\eta_i + 1)}{\partial \eta_i} = \frac{\partial A(\gamma_i)}{\partial \gamma_i} \times 1 = \\
&= \frac{\partial}{\partial \gamma_i} \sum_{i=1}^K \log \Gamma(\gamma_i) - \frac{\partial}{\partial \gamma_i} \log \Gamma\left(\sum_{j=1}^K \gamma_j\right) = \\
&= \Psi(\gamma_i) - \Psi\left(\sum_{j=1}^K \gamma_j\right),
\end{aligned} \tag{45}$$

dove abbiamo utilizzato la **funzione Digamma**, ossia la derivata logaritmica della funzione Gamma (facilmente calcolabile attraverso sviluppi in serie):

$$\Psi(z) = \frac{d}{dz} \log \Gamma(z) = \frac{\Gamma'(z)}{\Gamma(z)}. \tag{46}$$

Siamo pertanto giunti finalmente all'espressione finale del primo termine dell'energia libera:

$$\begin{aligned}
E_q \log[p(\theta|\alpha)] &= \log\left(\sum_{j=1}^K \alpha_j\right) - \\
&\quad - \sum_{i=1}^K \log \Gamma(\alpha_i) + \\
&\quad + \sum_{i=1}^K (\alpha_i - 1) \left[\Psi(\gamma_i) - \Psi\left(\sum_{j=1}^K \gamma_j\right) \right].
\end{aligned} \tag{47}$$

2. $E_q \log[p(\mathbf{z}|\theta)]$.

Per il secondo termine, la distribuzione delle variabili latenti è già stata esplicitata basandoci sul fatto che si tratta di una Multinomiale K -dimensionale su una singola prova, si veda la (13), e quindi:

$$p(\mathbf{z}|\theta) = \prod_{n=1}^N p(z_n|\theta) = \prod_{n=1}^N \prod_{i=1}^K \theta_i^{z_n^i}. \quad (48)$$

Ma allora:

$$\log p(\mathbf{z}|\theta) = \sum_{n=1}^N \sum_{i=1}^K z_n^i \log \theta_i, \quad (49)$$

e quindi, sfruttando l'indipendenza a priori tra le componenti della distribuzione variazionale:

$$\begin{aligned} E_q[\log p(\mathbf{z}|\theta)] &= E_q \left[\sum_{n=1}^N \sum_{i=1}^K z_n^i \log \theta_i \right] = \\ &= \sum_{n=1}^N \sum_{i=1}^K E_q[z_n^i \log \theta_i] = \\ &= \sum_{n=1}^N \sum_{i=1}^K E_q(z_n^i) E_q(\log \theta_i). \end{aligned} \quad (50)$$

La forma di $E_q(\log \theta_i)$ è stata esplicitata nella (45); invece, poiché $z_n|\phi$ ha distribuzione Multinomiale (su una singola prova) di parametri $(\phi_{n1}, \phi_{n2}, \dots, \phi_{nK})$, per la componente marginale z_n^i si ha immediatamente che:

$$E_q(z_n^i) = \phi_{ni}. \quad (51)$$

Pertanto:

$$E_q \log [p(\mathbf{z}|\theta)] = \sum_{n=1}^N \sum_{i=1}^K \phi_{ni} \left[\Psi(\gamma_i) - \Psi \left(\sum_{j=1}^K \gamma_j \right) \right]. \quad (52)$$

3. $E_q \log [p(\mathbf{w}|\mathbf{z}, \beta)]$.

Anche in questo caso abbiamo già esplicitato la distribuzione coinvolta nell'espressione (12), e quindi:

$$\begin{aligned}
E_q \log[p(\mathbf{w}|\mathbf{z}, \beta)] &= E_q \left[\sum_{n=1}^N \sum_{j=1}^{|\mathcal{V}|} w_n^j \log \beta_{ij} \right] = \\
&= \sum_{n=1}^N \sum_{j=1}^{|\mathcal{V}|} E_q [w_n^j \log \beta_{ij}] = \\
&= \sum_{n=1}^N \sum_{j=1}^{|\mathcal{V}|} w_n^j E_q [\log \beta_{ij}],
\end{aligned} \tag{53}$$

poiché w_n^j è ovviamente costante rispetto alla distribuzione variazionale. Si osservi ora che in $E_q[\log \beta_{ij}]$ la quantità aleatoria è l'indice i che determina il topic associato al token n -esimo. Ma poiché sotto la distribuzione variazionale q l'indicatore K -dimensionale z_n è una multinomiale con parametri $(\phi_{n1}, \phi_{n2}, \dots, \phi_{nK})$, ne segue che:

$$E_q[\log \beta_{ij}] = \sum_{i=1}^K \log \beta_{ij} \phi_{ni}, \tag{54}$$

e quindi:

$$E_q \log [p(\mathbf{w}|\mathbf{z}, \beta)] = \sum_{n=1}^N \sum_{j=1}^{|\mathcal{V}|} \sum_{i=1}^K w_n^j \phi_{ni} \log \beta_{ij}. \tag{55}$$

4. $E_q \log[q(\theta|\gamma)]$

Per questo termine, poiché $\theta|\gamma \stackrel{q}{\sim} \text{Dirichlet}_K(\gamma)$, i passaggi sono identici a quelli già visti per $E_q \log[p(\theta|\alpha)]$, cambiando ciò che va cambiato. Otteniamo, pertanto:

$$\begin{aligned}
E_q \log[q(\theta|\gamma)] &= \log \left(\sum_{j=1}^K \gamma_j \right) - \\
&\quad - \sum_{i=1}^K \log \Gamma(\gamma_i) + \\
&\quad + \sum_{i=1}^K (\gamma_i - 1) \left[\Psi(\gamma_i) - \Psi \left(\sum_{j=1}^K \gamma_j \right) \right].
\end{aligned} \tag{56}$$

5. $E_q \log[q(\mathbf{z}|\boldsymbol{\phi})]$

I calcoli relativi all'ultimo termine sono agevoli alla luce di quanto abbiamo già visto:

$$\begin{aligned}
E_q \log[q(\mathbf{z}|\boldsymbol{\phi})] &= \\
&= E_q \left[\log \prod_{n=1}^N q(z_n|\boldsymbol{\phi}_n) \right] = E_q \left[\log \prod_{n=1}^N \prod_{i=1}^K \phi_{ni} z_n^i \right] = \\
&= \sum_{n=1}^N \sum_{i=1}^K E_q[z_n^i] \log \phi_{ni} = \sum_{n=1}^N \sum_{i=1}^K \phi_{ni} \log \phi_{ni}.
\end{aligned} \tag{57}$$

Possiamo finalmente esplicitare l'espressione complessiva dell'energia libera (36), che dipende tanto dai parametri variazionali quanto dagli iperparametri α e β (oltre che, ovviamente, dai dati \mathbf{w}):

$$\begin{aligned}
\mathcal{L}(\boldsymbol{\gamma}, \boldsymbol{\phi}|\boldsymbol{\alpha}, \boldsymbol{\beta}) &= \log \Gamma \left(\sum_{j=1}^K \alpha_j \right) - \sum_{i=1}^K \log \Gamma(\alpha_i) \\
&+ \sum_{i=1}^K (\alpha_i - 1) \left[\Psi(\gamma_i) - \Psi \left(\sum_{j=1}^K \gamma_j \right) \right] + \\
&+ \sum_{n=1}^N \sum_{i=1}^K \phi_{ni} \left[\Psi(\gamma_i) - \Psi \left(\sum_{j=1}^K \gamma_j \right) \right] + \\
&+ \sum_{n=1}^N \sum_{j=1}^{|\mathcal{V}|} \sum_{i=1}^K w_n^j \phi_{ni} \log \beta_{ij} - \\
&- \log \Gamma \left(\sum_{j=1}^K \gamma_j \right) + \sum_{i=1}^K \log \Gamma(\gamma_i) - \\
&- \sum_{i=1}^K (\gamma_i - 1) \left[\Psi(\gamma_i) - \Psi \left(\sum_{j=1}^K \gamma_j \right) \right] - \\
&- \sum_{n=1}^N \sum_{i=1}^K \phi_{ni} \log \phi_{ni}.
\end{aligned} \tag{58}$$

5. L'algoritmo EM variazionale

Nell'espressione (58) dell'energia abbiamo potuto sopprimere l'indice relativo al documento d poiché i parametri variazionali sono locali: stante l'ipotesi di indipendenza a priori tra le componenti della distribuzione variazionale, l'energia libera complessiva per l'intero corpus può essere scritta come:

$$\sum_{d=1}^M \mathcal{L}(\gamma_d, \phi_d | \alpha, \beta), \quad (59)$$

dove ciascun documento ha i suoi propri specifici parametri variazionali (γ_d, ϕ_d) .

Come abbiamo già detto, dobbiamo massimizzare l'energia libera rispetto ai parametri variazionali per ottenere un minorante della log-verosimiglianza (non esprimibile) che sia il più vicino possibile ad essa, e dobbiamo poi utilizzare tale estremo inferiore come surrogato della log-verosimiglianza per ottenere le stime EB degli iperparametri. In altre parole, dobbiamo massimizzare l'energia libera simultaneamente tanto rispetto ai parametri variazionali quanto rispetto agli iperparametri.

Dal punto di vista numerico possiamo raggiungere agevolmente il nostro obiettivo attraverso il seguente **algoritmo EM variazionale**, che cicla iterativamente attraverso lo step di massimizzazione di (γ_d, ϕ_d) e quello di (α, β) :

- **E-step.** Per ciascun documento, ricerchiamo i valori ottimali (γ_d, ϕ_d) dei parametri variazionali, dati gli iperparametri (α, β) .
- **M-step.** Massimizziamo l'energia libera risultante (dati i parametri variazionali ottimali ottenuti allo step precedente) rispetto ai parametri del modello LDA (α, β) .

I due step vengono ripetuti iterativamente finché l'energia libera complessiva sull'intero corpus non si stabilizza ad un valore costante. Si noti che il primo step prende il nome di E-step poiché l'energia libera è un valore atteso rispetto alla distribuzione variazionale. Il secondo step è invece uno step di massimizzazione vera e propria dell'estremo inferiore della log-verosimiglianza marginale, poiché mira ad ottenere delle stime EB approssimate (approssimate in quanto la log-verosimiglianza marginale è approssimata dall'energia libera) dei parametri del modello.

Esplicitiamo ora l'E-step poiché i parametri variazionali sono locali possiamo sopprimere l'indicatore del documento d . Partiamo dalla ricerca dei valori ottimali di ϕ , che compare nei termini (52), (55) e (57) dell'energia libera. Osserviamo,

innanzitutto, che per ciascun token abbiamo un problema di massimo vincolato poiché deve essere soddisfatto l'ovvio vincolo:

$$\sum_{i=1}^K \phi_{ni} = 1, \quad \text{per ogni } n = 1, 2, \dots, N. \quad (60)$$

Poiché dobbiamo uguagliare a zero le derivate prime termine per termine, isoliamo i termini dell'energia libera che contengono il solo ϕ_{ni} , ed introduciamo il moltiplicatore di Lagrange per tenere conto del vincolo, ottenendo:

$$\begin{aligned} \mathcal{L}_{[\phi_{ni}]} &= \phi_{ni} \left[\Psi(\gamma_i) - \Psi \left(\sum_{j=1}^K \gamma_j \right) \right] + \\ &+ \sum_{j=1}^{|\mathcal{V}|} w_n^j \phi_{ni} \log \beta_{ij} - \\ &- \phi_{ni} \log \phi_{ni} - \lambda \left[\sum_{j=1}^K \phi_{nj} - 1 \right]. \end{aligned} \quad (61)$$

Osservando che per l' n -esimo token esiste un unico indice v per il quale $w_n^v = 1$ e $w_n^u = 0$ per $u \neq v$, possiamo scrivere:

$$\begin{aligned} \mathcal{L}_{[\phi_{ni}]} &= \phi_{ni} \left[\Psi(\gamma_i) - \Psi \left(\sum_{j=1}^K \gamma_j \right) \right] + \phi_{ni} \log \beta_{iv} - \\ &- \phi_{ni} \log \phi_{ni} - \lambda \left[\sum_{j=1}^K \phi_{nj} - 1 \right]. \end{aligned} \quad (62)$$

Possiamo ora facilmente derivare la Lagrangiana, ottenendo:

$$\begin{aligned} \frac{\partial \mathcal{L}_{[\phi_{ni}]} }{\partial \phi_{ni}} &= \left[\Psi(\gamma_i) - \Psi \left(\sum_{j=1}^K \gamma_j \right) \right] + \\ &+ \log \beta_{iv} - \log \phi_{ni} - \phi_{ni} \frac{1}{\phi_{ni}} - \lambda. \end{aligned} \quad (63)$$

Uguagliando a zero la derivata parziale prima abbiamo:

$$\log \phi_{ni} = -1 - \lambda + \log \beta_{iv} + \left[\Psi(\gamma_i) - \Psi\left(\sum_{j=1}^K \gamma_j\right) \right] \quad (64)$$

ossia, ancora:

$$\phi_{ni} = \exp(-1 - \lambda) \beta_{iv} \exp \left[\Psi(\gamma_i) - \Psi\left(\sum_{j=1}^k \gamma_j\right) \right], \quad (65)$$

e quindi:

$$\phi_{ni} \propto \beta_{iv} \exp \left[\Psi(\gamma_i) - \Psi\left(\sum_{j=1}^k \gamma_j\right) \right], \quad (66)$$

Il moltiplicatore di Lagrange scompare automaticamente, poiché basta normalizzare ad 1 per tenere conto del vincolo (58):

$$\phi_{ni} \leftarrow \frac{\phi_{ni}}{\sum_{j=1}^K \phi_{nj}}, \quad (67)$$

e quindi il termine $\exp(-1 - \lambda)$ compare tanto al numeratore quanto al denominatore, e può essere semplificato in quanto costante (rispetto all'indice i). Si osservi che possiamo ulteriormente semplificare nel modo seguente:

$$\begin{aligned} \phi_{ni} &\propto \beta_{iv} \exp \left[\Psi(\gamma_i) - \Psi\left(\sum_{j=1}^K \gamma_j\right) \right] = \\ &= \beta_{iv} \exp \Psi(\gamma_i) \exp \left[-\Psi\left(\sum_{j=1}^K \gamma_j\right) \right] \propto \\ &\propto \beta_{iv} \exp \Psi(\gamma_i), \end{aligned} \quad (68)$$

poiché l'ultimo fattore non dipende dall'indice i , e quindi anch'esso viene eliminato durante la normalizzazione ad 1.

Scriviamo ora i termini dell'energia libera che includono γ_i . In questo caso non abbiamo bisogno di introdurre un moltiplicatore di Lagrange poiché non abbiamo alcun vincolo da soddisfare. Abbiamo, pertanto:

$$\begin{aligned}
\mathcal{L}_{[\gamma_i]} &= (\alpha_i - 1) \left[\Psi(\gamma_i) - \Psi\left(\sum_{j=1}^K \gamma_j\right) \right] + \\
&+ \sum_{n=1}^N \phi_{ni} \left[\Psi(\gamma_i) - \Psi\left(\sum_{j=1}^K \gamma_j\right) \right] - \\
&- \log \Gamma\left(\sum_{j=1}^K \gamma_j\right) + \log \Gamma(\gamma_i) - \\
&- (\gamma_i - 1) \left[\Psi(\gamma_i) - \Psi\left(\sum_{j=1}^K \gamma_j\right) \right] = \\
&= \left[\Psi(\gamma_i) - \Psi\left(\sum_{j=1}^K \gamma_j\right) \right] \left(\alpha_i - 1 + \sum_{n=1}^N \phi_{ni} - \gamma_i + 1 \right) - \\
&- \log \Gamma\left(\sum_{j=1}^K \gamma_j\right) + \log \Gamma(\gamma_i) = \\
&= \left[\Psi(\gamma_i) - \Psi\left(\sum_{j=1}^K \gamma_j\right) \right] \left(\alpha_i + \sum_{n=1}^N \phi_{ni} - \gamma_i \right) - \\
&- \log \Gamma\left(\sum_{j=1}^K \gamma_j\right) + \log \Gamma(\gamma_i).
\end{aligned} \tag{69}$$

Prendendo le derivate parziali rispetto a γ_i ed uguagliandole a 0, abbiamo:

$$\begin{aligned}
\frac{\partial \mathcal{L}_{[\gamma_i]}}{\partial \gamma_i} &= \Psi'(\gamma_i) \left(\alpha_i + \sum_{n=1}^N \phi_{ni} - \gamma_i \right) + \Psi(\gamma_i)(-1) - \\
&- \Psi' \left(\sum_{j=1}^K \gamma_j \right) \left(\alpha_i + \sum_{n=1}^N \phi_{ni} - \gamma_i \right) - \Psi \left(\sum_{j=1}^K \gamma_j \right) (-1) - \\
&- \Psi \left(\sum_{j=1}^K \gamma_j \right) + \Psi(\gamma_i) = \tag{70} \\
&= \Psi'(\gamma_i) \left(\alpha_i + \sum_{n=1}^N \phi_{ni} - \gamma_i \right) - \\
&- \Psi' \left(\sum_{j=1}^K \gamma_j \right) \left(\alpha_i + \sum_{n=1}^N \phi_{ni} - \gamma_i \right) = 0,
\end{aligned}$$

ossia:

$$\begin{aligned}
\Psi'(\gamma_i) \left(\alpha_i + \sum_{n=1}^N \phi_{ni} - \gamma_i \right) &= \\
&= \Psi' \left(\sum_{j=1}^K \gamma_j \right) \left(\alpha_i + \sum_{n=1}^N \phi_{ni} - \gamma_i \right), \tag{71}
\end{aligned}$$

il che è possibile se e solo se:

$$\alpha_i + \sum_{n=1}^N \phi_{ni} - \gamma_i = 0, \tag{72}$$

che fornisce immediatamente l'equazione di aggiornamento per γ_i :

$$\gamma_i = \alpha_i + \sum_{n=1}^N \phi_{ni}. \tag{73}$$

Ovviamente, questa equazione di aggiornamento va accoppiata con la precedente su ϕ_{ni} . Anche in questo caso dobbiamo pertanto ciclare tra l'aggiornamento su ϕ_{ni}

(per $n = 1, 2, \dots, N$ e $i = 1, 2, \dots, K$) e quello su γ_i (per $i = 1, 2, \dots, K$), finché il valore dell'energia libera per lo step corrente non si stabilizza. Lo pseudo-codice dell'E-step per un singolo documento riportato qui sotto è quello originariamente proposto in Blei et al., (2003):

Figura 5. Pseudo-codice dell'E-step nell'algoritmo EM variazionale utilizzato per la stima dei parametri del modello LDA.

```

1  inizializza  $\phi_{ni}^{(0)} \leftarrow \frac{1}{K}$  per ogni  $i$  e  $n$ 
2  inizializza  $\gamma_i^{(0)} \leftarrow \alpha_i + \frac{N}{K}$  per ogni  $i$ 
3  repeat
4    for  $n = 1$  to  $N$  do
5      for  $i = 1$  to  $K$  do
6         $\phi_{ni}^{(t+1)} \leftarrow \beta_{iw_n} \exp(\Psi(\gamma_i^{(t)}))$ 
7      end
8      normalizza la somma dei  $\phi_{ni}^{(t+1)}$  a 1
9    end
10    $\gamma^{(t+1)} \leftarrow \alpha + \sum_{n=1}^N \phi_n^{(t+1)}$ 
11  until convergence

```

Alcune osservazioni degne di nota sono le seguenti:

- per ciascun token, gli iperparametri della distribuzione Multinomiale variazionale sui topic sono usualmente inizializzati in modo uniforme rispetto al numero di topic;
- gli iperparametri della distribuzione di Dirichlet variazionale delle proporzioni dei topic sono inizializzati come $\alpha_i + N/K$: ciò corrisponde ad aggiornare il valore corrente dell'iperparametro globale α_i in base al numero atteso di token per ciascun topic sotto una distribuzione Multinomiale con probabilità uniformi rispetto al numero dei topic;
- le inizializzazioni scelte sono ovviamente arbitrarie; inoltre, è ovviamente ridondante inizializzare entrambi i parametri variazionali. Poiché nell'E-step l'algoritmo cicla tra questi due insiemi di parametri, è sufficiente inizializzarne solo uno;
- con abuso di notazione, nello step 6 abbiamo scritto β_{iw_n} , identificando w_n con l'unico indice v per il quale $w_n^v = 1$;
- l'aggiornamento di γ allo step 10 è ovviamente scritto in forma vettoriale, e deve essere sviluppato componente per componente.

Vediamo ora i dettagli dell'M-step. Come abbiamo già ampiamente osservato in precedenza, i parametri α e β non sono locali, e quindi è necessario considerare

anche l'indice d del documento, e sommare tutti i termini dell'energia libera documento per documento. Se consideriamo β , l'unico termine dell'energia nel quale compare β è il terzo. Abbiamo però dei vincoli derivanti dal fatto che ogni riga della matrice β è una distribuzione di probabilità sul vocabolario dei termini, ossia:

$$\sum_{j=1}^{|V|} \beta_{ij} = 1. \quad (74)$$

Pertanto, isolando solo il termine che contiene β_{ij} abbiamo la seguente lagrangiana:

$$\begin{aligned} \mathcal{L}_{[\beta_{ij}]} &= \sum_{d=1}^M \sum_{n=1}^{N_d} \sum_{j=1}^{|V|} \sum_{i=1}^K w_{dn}^j \phi_{dni} \log \beta_{ij} - \\ &- \sum_{i=1}^K \lambda_i \left[\sum_{j=1}^{|V|} \beta_{ij} - 1 \right]. \end{aligned} \quad (75)$$

Prendiamo le derivate parziali; per il primo addendo abbiamo:

$$\begin{aligned} \frac{\partial}{\partial \beta_{ij}} \sum_{d=1}^M \sum_{n=1}^{N_d} \sum_{s=1}^{|V|} \sum_{t=1}^K w_{dn}^s \phi_{dnt} \log \beta_{ts} &= \\ = \sum_{d=1}^M \sum_{n=1}^{N_d} w_{dn}^j \phi_{dni} \frac{1}{\beta_{ij}} \frac{\partial \beta_{ij}}{\partial \beta_{ij}} &= \\ = \sum_{d=1}^M \sum_{n=1}^{N_d} w_{dn}^j \phi_{dni} \frac{1}{\beta_{ij}}. \end{aligned} \quad (76)$$

Per il termine contenente i moltiplicatori di Lagrange abbiamo, invece:

$$\begin{aligned} \frac{\partial}{\partial \beta_{ij}} \sum_{t=1}^K \lambda_t \left[\sum_{s=1}^{|V|} \beta_{ts} - 1 \right] &= \\ = \frac{\partial}{\partial \beta_{ij}} \sum_{t=1}^K \sum_{s=1}^{|V|} \lambda_t \beta_{st} - \frac{\partial}{\partial \beta_{ij}} \sum_{t=1}^K \lambda_t &= \lambda_i - 0 = \lambda_i. \end{aligned} \quad (77)$$

Pertanto, l'equazione:

$$\frac{\partial \mathcal{L}_{[\beta_{ij}]}}{\partial \beta_{ij}} = 0, \quad (78)$$

diventa:

$$\sum_{d=1}^M \sum_{n=1}^{N_d} w_{dn}^j \phi_{dni} \frac{1}{\beta_{ij}} = -\lambda_i. \quad (79)$$

Poiché β_{ij} non dipende dagli indici delle sommatorie, possiamo scrivere:

$$\frac{1}{\beta_{ij}} = -\frac{\lambda_i}{\sum_{d=1}^M \sum_{n=1}^{N_d} w_{dn}^j \phi_{dni} \frac{1}{\beta_{ij}}}, \quad (80)$$

e quindi, infine:

$$\beta_{ij} = -\frac{1}{\lambda_i} \sum_{d=1}^M \sum_{n=1}^{N_d} w_{dn}^j \phi_{dni}. \quad (81)$$

Anche in questo caso i moltiplicatori di Lagrange si eliminano tenendo conto del fatto che sommando rispetto all'indice j (per i prefissato) le quantità β_{ij} sommano ad 1, essendo distribuzioni di probabilità, e quindi la normalizzazione avviene, esattamente come abbiamo fatto con ϕ_n , ponendo:

$$\beta_{ij} \propto \sum_{d=1}^M \sum_{n=1}^{N_d} w_{dn}^j \phi_{dni}, \quad (82)$$

e normalizzando nel modo seguente:

$$\beta_{ij} \leftarrow \frac{\beta_{ij}}{\sum_{j=1}^{|V|} \beta_{ij}}. \quad (83)$$

Il fattore $-\lambda_i^{-1}$, che compare tanto al numeratore quanto al denominatore, viene semplificato.

Per completare la trattazione, isoliamo dall'energia libera i termini che contengono α . Abbiamo:

$$\begin{aligned} \mathcal{L}_{[\alpha]} &= \sum_{d=1}^M \left\{ \log \Gamma \left(\sum_{j=1}^K \alpha_j \right) - \sum_{i=1}^K \log \Gamma(\alpha_t) + \right. \\ &\quad \left. + \sum_{t=1}^K (\alpha_t - 1) \left[\Psi(\gamma_{at}) - \Psi \left(\sum_{j=1}^K \gamma_{dj} \right) \right] \right\}, \end{aligned} \quad (84)$$

e pertanto:

$$\begin{aligned} \frac{\partial}{\partial \alpha_i} \mathcal{L}_{[\alpha]} &= M \left[\Psi \left(\sum_{j=1}^K \alpha_j \right) - \Psi(\alpha_i) \right] + \\ &\quad + \sum_{d=1}^M \frac{\partial}{\partial \alpha_i} \sum_{t=1}^K \alpha_t \left[\Psi(\gamma_{at}) - \Psi \left(\sum_{j=1}^K \gamma_{dj} \right) \right] = \\ &= M \left[\Psi \left(\sum_{j=1}^K \alpha_j \right) - \Psi(\alpha_i) \right] + \\ &\quad + \sum_{d=1}^M \left[\Psi(\gamma_{at}) - \Psi \left(\sum_{j=1}^K \gamma_{dj} \right) \right]. \end{aligned} \quad (85)$$

Uguagliando a zero questa derivate parziale, abbiamo una equazione in α_i da risolvere, che è non lineare e la relativa soluzione non è direttamente esplicitabile.

Per risolvere questa equazione possiamo utilizzare una versione $\mathcal{O}(K)$ in tempo lineare dell'algoritmo di Newton-Raphson (NR), che nella sua formulazione standard è di ordine $\mathcal{O}(K^3)$ a causa delle inversioni matriciali che sono richieste. Ricordiamo che nella sua forma generale il metodo di NR permette di risolvere l'equazione vettoriale $g(\varphi) = 0$, con $\varphi \in \mathbb{R}^K$. L'equazione di aggiornamento del metodo è la seguente:

$$\varphi^{(s+1)} = \varphi^{(s)} - \{H(\varphi)\}^{-1} u(\varphi^{(s)}), \quad (86)$$

dove $H(\varphi)$ è la **matrice Hessiana**, la cui inversione è responsabile dell'elevata complessità computazionale dell'algoritmo, mentre $u(\varphi)$ è il **gradiente**:

$$u(\varphi) = \begin{pmatrix} \frac{\partial}{\partial \varphi_1} g(\varphi) \\ \frac{\partial}{\partial \varphi_2} g(\varphi) \\ \vdots \\ \frac{\partial}{\partial \varphi_K} g(\varphi) \end{pmatrix}. \quad (87)$$

Facciamo l'ipotesi che l'Hessiano possa essere scritto nel modo seguente (per semplicità sopprimiamo la dipendenza dall'argomento φ):

$$H = \text{diag}(h) + \mathbb{1}z\mathbb{1}^\top, \quad (88)$$

con $h \in \mathbb{R}^K$, $\mathbb{1} \in \mathbb{R}^K$ è il vettore K -dimensionale i cui elementi sono tutti 1, e z è uno scalare. Se questa ipotesi è vera, la struttura della matrice Hessiana è evidentemente del tipo:

$$H = A + UCV, \quad (89)$$

dove:

$$\begin{aligned} A &\equiv \text{diag}(h) \in \mathbb{R}^{K \times K} \\ U &\equiv \mathbb{1} \in \mathbb{R}^{K \times 1} \\ C &\equiv z \in \mathbb{R}^{1 \times 1} \\ V &\equiv \mathbb{1}^\top \in \mathbb{R}^{1 \times K}. \end{aligned} \quad (90)$$

Possiamo allora applicare la ben nota **identità di Woodbury** (matrix inversion lemma), grazie alla quale l'inversa di una matrice H che ha la struttura che abbiamo indicata nella (89) è la seguente:

$$H^{-1} = A^{-1} - A^{-1}U(C^{-1} + VA^{-1}U)^{-1}VA^{-1}, \quad (91)$$

ossia, sostituendo ciò che va sostituito:

$$\begin{aligned} H^{-1} &= \text{diag}(h)^{-1} - \text{diag}(h)^{-1}\mathbb{1} \times \\ &\times \left[z^{-1} + \mathbb{1}^\top \text{diag}(h)^{-1}\mathbb{1} \right] \mathbb{1}^\top \text{diag}(h)^{-1}. \end{aligned} \quad (92)$$

Ma osserviamo che:

$$\begin{aligned}
 & \mathbf{1}^\top \text{diag}(h)^{-1} \mathbf{1} = \\
 & = (1 \quad 1 \quad \dots \quad 1) \begin{pmatrix} \frac{1}{h_1} & 0 & \dots & 0 \\ 0 & \frac{1}{h_2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{1}{h_K} \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} = \\
 & = \left(\frac{1}{h_1} \quad \frac{1}{h_2} \quad \dots \quad \frac{1}{h_K} \right) \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} = \\
 & = \sum_{j=1}^K h_j^{-1},
 \end{aligned} \tag{93}$$

e pertanto:

$$\begin{aligned}
 H^{-1} & = \text{diag}(h)^{-1} - \text{diag}(h)^{-1} \mathbf{1} \mathbf{1}^\top \text{diag}(h)^{-1} \times \\
 & \times \frac{1}{z^{-1} + \sum_{j=1}^K h_j^{-1}}.
 \end{aligned} \tag{94}$$

Sviluppando l'altro prodotto matriciale, con semplici calcoli otteniamo:

$$\begin{aligned}
& \text{diag}(h)^{-1} \mathbb{1} \mathbb{1}^T \text{diag}(h)^{-1} = \\
& = \begin{pmatrix} \frac{1}{h_1} & 0 & \dots & 0 \\ 0 & \frac{1}{h_2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{1}{h_K} \end{pmatrix} \begin{pmatrix} 1 & 1 & \dots & 1 \\ 1 & 1 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \end{pmatrix} \begin{pmatrix} \frac{1}{h_1} & 0 & \dots & 0 \\ 0 & \frac{1}{h_2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{1}{h_K} \end{pmatrix} = \\
& = \begin{pmatrix} \frac{1}{h_1} & \frac{1}{h_1} & \dots & \frac{1}{h_1} \\ \frac{1}{h_2} & \frac{1}{h_2} & \dots & \frac{1}{h_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{1}{h_K} & \frac{1}{h_K} & \dots & \frac{1}{h_K} \end{pmatrix} \begin{pmatrix} \frac{1}{h_1} & 0 & \dots & 0 \\ 0 & \frac{1}{h_2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{1}{h_K} \end{pmatrix} = \\
& = \begin{pmatrix} \frac{1}{h_1^2} & \frac{1}{h_1 h_2} & \dots & \frac{1}{h_1 h_K} \\ \frac{1}{h_2 h_1} & \frac{1}{h_2^2} & \dots & \frac{1}{h_2 h_K} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{1}{h_K h_1} & \frac{1}{h_K h_2} & \dots & \frac{1}{h_K^2} \end{pmatrix}.
\end{aligned} \tag{95}$$

Pertanto, quando nell'equazione di aggiornamento del metodo di NR andiamo a moltiplicare H^{-1} per il gradiente \mathbf{u} , la i -esima componente del vettore risultante sarà la seguente:

$$\begin{aligned}
(H^{-1} \mathbf{u})_i &= \frac{u_i}{h_i} - \frac{\sum_{j=1}^K u_j / h_i h_j}{z^{-1} + \sum_{j=1}^K h_j^{-1}} = \\
&= \frac{u_i}{h_i} - \frac{\sum_{j=1}^K u_j / h_j}{z^{-1} + \sum_{j=1}^K h_j^{-1}} = \\
&= \frac{u_i}{h_i} + \frac{c}{h_i} = \frac{g_i - c}{h_i},
\end{aligned} \tag{96}$$

dove:

$$c = \frac{\sum_{j=1}^K u_j / h_j}{z^{-1} + \sum_{j=1}^K h_j^{-1}}, \quad (97)$$

è una quantità costante che va calcolata una volta per tutte. Ma allora, il calcolo di ciascun elemento del prodotto $H^{-1}u$ ha bisogno solo del set di $2K$ valori u_i ed h_i , da cui la complessità in tempo lineare di questa versione specializzata dell'algoritmo di Newton-Rapshon.

Osserviamo ora che:

$$\begin{aligned} \frac{\partial \mathcal{L}_{[\alpha]}}{\partial \alpha_i^2} &= M\Psi' \left(\sum_{j=1}^K \alpha_j \right) - M\Psi'(\alpha_i), \\ \frac{\partial \mathcal{L}_{[\alpha]}}{\partial \alpha_i \alpha_j} &= M\Psi' \left(\sum_{j=1}^K \alpha_j \right), \quad i \neq j, \end{aligned} \quad (98)$$

e, di conseguenza, il generico elemento della matrice Hessiana ha espressione:

$$H_{ij} = -\delta_{ij} M\Psi'(\alpha_i) + M\Psi' \left(\sum_{j=1}^K \alpha_j \right), \quad (99)$$

che corrisponde esattamente alla forma richiesta per H affinché si possa utilizzare la versione in tempo lineare dell'algoritmo di NR, prendendo:

$$\begin{aligned} h_i &= -M\Psi'(\alpha_i), \\ z &= M\Psi' \left(\sum_{j=1}^K \alpha_j \right). \end{aligned} \quad (100)$$

Naturalmente, una volta che abbiamo determinato i parametri variazionali ottimali (ϕ^*, γ^*) , possiamo utilizzare direttamente la distribuzione variazionale $q(\theta, \mathbf{z} | \phi^*, \gamma^*)$, che è la miglior approssimazione possibile della distribuzione a posteriori, per ottenere le stime della quantità di interesse. Sotto la distribuzione variazionale questo compito è particolarmente agevole, in quanto le componenti della distribuzione variazionale si fattorizzano. Ad esempio, poiché $q(\theta_d | \gamma_d^*)$ è una

distribuzione di Dirichlet, la stima delle proporzioni di ciascun topic nel documento d -esimo (**per-document topic proportions**) può essere ottenuta come:

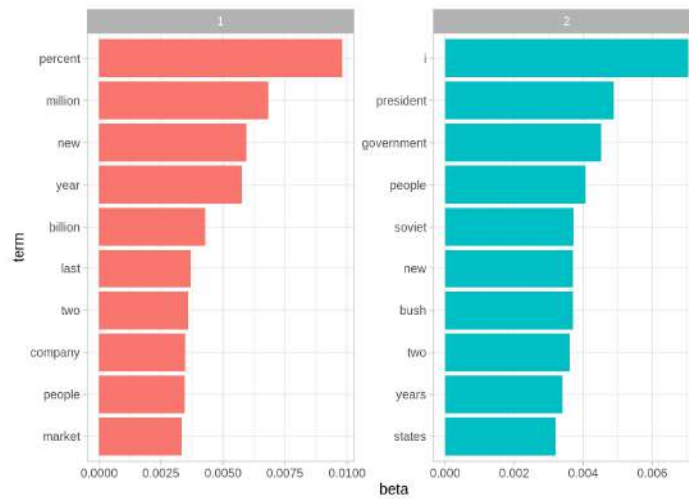
$$\hat{\theta}_{di} = E_q(\theta_{di}|\gamma^*) = \frac{\gamma_{di}^*}{\sum_{j=1}^K \gamma_{dj}^*}. \quad i = 1, 2, \dots, K, \quad (101)$$

così come le stime $\hat{\beta}_{ij} = \beta_{ij}^*$ del parametro globale β ci forniscono le stime EB delle distribuzioni di probabilità sui topic (**per-term topic probabilities**). Possiamo anche ottenere facilmente la stima della probabilità che il token n -esimo nel documento d sia generato dall' i -esimo topic (**per-word topic assignment proportions**):

$$\hat{z}_{dni} = E_q(z_{dni}|\phi^*) = \phi_{dni}^*. \quad (102)$$

Queste quantità, sebbene fondamentali per la struttura generativa del modello, sono di minore rilevanza a fini pratici, poiché è quasi inutile conoscere la stima delle probabilità multinomiali degli indicatori latenti, così come è inutile conoscere le assegnazioni ai topic per ciascun termine di ciascun documento (avremmo troppa informazione difficile da sintetizzare). Invece, un riassunto della semantica del corpus a nostra disposizione può essere ottenuto attraverso le stime $\hat{\beta}_{ij}$, ad esempio plot-tando per ciascun topic i termini con le probabilità più elevate. Nella Figura 6 riportiamo, a titolo di esempio, i dieci top-token di un caso di studio (ossia i dieci token che hanno le per-term topic probabilities più elevate) per due topic:

Figura 6. Esempio relativo alla probabilità di occorrenza dei dieci top-token, ossia i dieci token che hanno le probabilità di occorrenza più elevate, per due topic.



Similmente, la distribuzione $(\hat{\theta}_{d1}, \hat{\theta}_{d2}, \dots, \hat{\theta}_{di})$ contiene, per il documento d , la proporzione dei termini che sono stati assegnati a ciascuno dei K topic. Poiché l'obiettivo della classificazione testuale è proprio quello di poter assegnare ogni documento all'argomento a cui appartiene con più probabilità, con il modello LDA questo obiettivo può essere raggiunto in modo non supervisionato assegnando il documento d a quel topic che ha generato la proporzione più alta di termini nel documento.

6. Inferenza variazionale nel modello full Bayes

Come abbiamo detto più volte, nel modello LDA considerato fino ad ora gli iperparametri α e β sono considerati fissi ma incogniti, e stimati dai dati attraverso l'approccio EB approssimato, ottenuto surrogando la verosimiglianza marginale con l'energia libera.

Naturalmente, possiamo considerare in alternativa un approccio **full Bayes**, nel quale ogni iperparametro ha una sua propria distribuzione di probabilità a priori completamente specificata, con iperparametri prefissati ad un valore ritenuto ragionevole sulla base dell'informazione disponibile, ovvero prefissati in modo tale che essi riflettano una situazione di ignoranza sulla quantità che stiamo modellando.

In questo modo, le distribuzioni sui topic sul vocabolario dei termini possono essere modellizzate attraverso K distribuzioni di Dirichlet indipendenti:

$$p(\beta_{1:K}|\eta) = \prod_{i=1}^K p(\beta_i|\eta), \quad (103)$$

con:

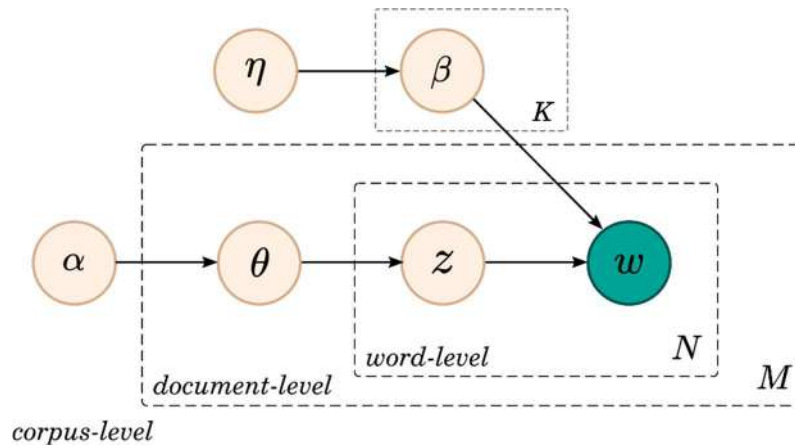
$$\beta_i \sim \text{Dirichlet}_{|V|}(\eta), \quad i = 1, 2, \dots, K, \quad (104)$$

In questa scrittura abbiamo utilizzato la notazione $\beta_{1:K}$ per sottolineare che ciascuna riga della matrice β viene dotata della propria distribuzione di probabilità a priori. L'iperparametro η indica un vettore costante $|V|$ -dimensionale, nel quale ciascun elemento è identicamente uguale allo scalare η , ed è comune a tutte le distribuzioni sui topic. In altre parole, stiamo assumendo che ciascuna distribuzione β_i sia modellizzata attraverso una distribuzione a priori di Dirichlet **simmetrica** che riflette a priori l'ignoranza sull'importanza assunta da ciascun token all'interno di un topic.

Nella specificazione del modello LDA full Bayes, di solito si postula una ipotesi di simmetria anche per l'altro iperparametro α , ponendo $\alpha_i = \alpha$ per $i = 1, 2, \dots, K$. In questo caso, una notazione vettoriale corretta per indicare tale iperparametro sarebbe $\mathbf{1}^\top \alpha$, così come $\mathbf{1}^\top \eta$ è la notazione vettoriale corretta per l'altro iperparametro: tuttavia, per semplicità continueremo a indicare tali quantità semplicemente con α ed η , avendo però ben chiaro che esse indicano vettori costanti. Supporremo, inoltre, che tanto α quanto η siano prefissati (ad esempio $\alpha = \eta = 1$). Infine, tutto quanto diremo per il modello full Bayes rimane formalmente valido anche nel caso in cui le distribuzioni a priori su $\beta_{1:K}$ e θ_i non siano Dirichlet simmetriche (con iperparametri completamente specificati). Naturalmente, in quest'ultimo caso avremo un impatto sull'inferenza a posteriori che viene esplorato, per esempio, in Wallach et al. (2009).

Graficamente, questa versione estesa del modello LDA può essere rappresentata mediante il seguente grafo orientato:

Figura 7. Il modello LDA nella versione Full Bayes.



Se vogliamo considerare un approccio variazionale all'inferenza, anche in questo caso utilizzeremo una distribuzione variazionale di tipo mean-field, nel quale tutte le componenti marginali si fattorizzano e sono indipendenti:

$$\begin{aligned}
 & q(\beta_{1:K}, \mathbf{z}, \theta | \lambda, \phi, \gamma) = \\
 & = \prod_{i=1}^K q(\beta_i | \lambda_i) \prod_{d=1}^M q(\theta_d, \mathbf{z}_d | \phi_d, \gamma_d), \tag{105}
 \end{aligned}$$

dove $q(\theta_d, \mathbf{z}_d | \phi_d, \gamma_d)$ è la stessa distribuzione variazionale che abbiamo visto in precedenza. Si noti che:

- abbiamo soppresso la dipendenza dagli iperparametri α e η , poiché abbiamo detto che essi sono prefissati;
- λ è un parametro variazionale globale (che non dipende da d), le cui componenti sono $\lambda_i \in \mathbb{R}^{|V|}$. Quindi, λ può essere visto come una matrice $K \times |V|$;
- come nel caso precedente, ϕ e γ sono parametri variazionali locali, che variano documento per documento.

Scriviamo ora direttamente l'energia libera totale valida per l'intero corpus, tenendo ovviamente presente che la parte dipendente da $\beta_{1:K}$ non va sommata su d (poiché questa distribuzione non dipende da d). Abbiamo:

$$\mathcal{L}(\gamma, \phi, \lambda) = E_q[\log p(\theta, \mathbf{z}, \beta, \mathbf{w})] - E_q[\log q(\theta, \mathbf{z}, \beta)], \quad (106)$$

e data la struttura di dipendenza condizionale delle distribuzioni coinvolte e sviluppando, otteniamo

$$\begin{aligned} \mathcal{L}(\gamma, \phi, \lambda) &= E_q[\log p(\beta_{1:K} | \eta)] + \\ &+ \sum_{d=1}^M \{E_q \log[p(\theta_d | \alpha)] + \\ &+ E_q \log[p(\mathbf{z}_d | \theta_d)] + \\ &+ E_q \log[p(\mathbf{w}_d | \mathbf{z}_d, \beta_{1:K})]\} - \\ &- E_q[\log q(\beta_{1:K} | \lambda)] - \\ &- \sum_{d=1}^M \{E_q \log[q(\theta_d | \gamma_d)] - \\ &- E_q \log[q(\mathbf{z}_d | \phi_d)]\}. \end{aligned} \quad (107)$$

Esplicitiamo uno per uno i termini coinvolti.

1. $E_q[\log p(\beta_{1:K}|\eta)]$

Innanzitutto, poiché ciascuna componente β_i ha distribuzione di Dirichlet simmetrica:

$$\begin{aligned} \log p(\beta_i|\eta) &= \log \frac{\Gamma(\sum_{j=1}^{|V|} \beta_{ij})}{\prod_{j=1}^{|V|} \Gamma(\beta_{ij})} \beta_{i1}^{(\eta-1)} \beta_{i2}^{(\eta-1)} \dots \beta_{i|V|}^{(\eta-1)} = \\ &= \log \Gamma(\eta|V) - |V| \log \Gamma(\eta) - \\ &\quad - (\eta - 1) \sum_{j=1}^{|V|} \log \beta_{ij}, \end{aligned} \quad (108)$$

e pertanto:

$$\begin{aligned} E_q[\log p(\beta_{1:K}|\eta)] &= E_q \left[\prod_{i=1}^K \log p(\beta_i|\eta) \right] = \\ &= E_q \left[\log \Gamma(\eta|V) - |V| \log \Gamma(\eta) \right. \\ &\quad \left. - \sum_{i=1}^K (\eta - 1) \sum_{j=1}^{|V|} \log \beta_{ij} \right] = \\ &= \log \Gamma(\eta|V) - |V| \log \Gamma(\eta) + \\ &\quad + \sum_{i=1}^K \sum_{j=1}^{|V|} (\eta - 1) E_q[\log \beta_{ij}] = \\ &= \log \Gamma(\eta|V) - |V| \log \Gamma(\eta) + \\ &\quad + \sum_{i=1}^K \sum_{j=1}^{|V|} (\eta - 1) \left[\Psi(\lambda_{ij}) - \Psi \left(\sum_{s=1}^{|V|} \lambda_{is} \right) \right]. \end{aligned} \quad (109)$$

2. $E_q \log[p(\theta_d|\alpha)]$

I calcoli sono identici a quelli già visti in precedenza, con $\alpha_i = \alpha$ per $i = 1, 2, \dots, K$, Dunque:

$$\begin{aligned} E_q \log [p(\theta_d|\alpha)] &= \log \Gamma(K\alpha) - K \log \Gamma(\alpha) + \\ &\quad + \sum_{i=1}^K (\alpha - 1) \left[\Psi(\gamma_{di}) - \Psi \left(\sum_{j=1}^K \gamma_{dj} \right) \right]. \end{aligned} \quad (110)$$

3. $\log[p(\mathbf{z}_d|\theta_d)]$

Anche in questo caso, come già visto in precedenza:

$$\begin{aligned} & E_q \log [p(\mathbf{z}_d|\theta_d)] = \\ & = \sum_{i=1}^K \sum_{n=1}^{N_d} \phi_{dni} \left[\Psi(\gamma_{di}) - \Psi \left(\sum_{j=1}^K \gamma_{dj} \right) \right]. \end{aligned} \quad (111)$$

Tuttavia, vogliamo mettere una formulazione alternativa basata sull'identità:

$$\sum_{n=1}^{N_d} \phi_{dni} = \sum_{j=1}^{|V|} n_{dj} \phi_{dji}, \quad (112)$$

dove n_{dj} è il numero di volte che il token j appare nel documento d . Dunque:

$$\begin{aligned} & E_q \log [p(\mathbf{z}_d|\theta_d)] = \\ & = \sum_{i=1}^K \sum_{j=1}^{|V|} n_{dj} \phi_{dji} \left[\Psi(\gamma_{di}) - \Psi \left(\sum_{j=1}^K \gamma_{dj} \right) \right], \end{aligned} \quad (113)$$

che ha l'ovvio vantaggio di esprimere il valore atteso in funzione delle frequenze assolute con le quali i token appaiono in ciascun documento.

4. $E_q \log[p(\mathbf{w}_d|\mathbf{z}_d, \beta_{1:K})]$

Con calcoli assolutamente identici a quelli già visti:

$$\begin{aligned} & E_q \log [p(\mathbf{w}_d|\mathbf{z}_d, \beta_{1:K})] = \\ & = \sum_{i=1}^K \sum_{j=1}^{|V|} \sum_{n=1}^{N_d} w_n^j \phi_{dni} E_q [\log \beta_{ij}]. \end{aligned} \quad (114)$$

Anche in questo caso vogliamo modificare questa espressione utilizzando un'altra ovvia identità, simile alla (112), e cioè:

$$\sum_{n=1}^{N_d} w_n^j \phi_{dni} = n_{dj} \phi_{dji}, \quad (115)$$

che sussiste in quanto nella sommatoria a sinistra j è prefissato, e quindi sommare w_n^j su tutti i termini del documento equivale a contare il numero di volte che il token j appare nel documento d . Inoltre, tenendo presente che sotto la distribuzione variazionale β_i ha distribuzione di Dirichlet con iperparametro variazionale λ_i :

$$\begin{aligned} & E_q \log [p(\mathbf{w}_d | \mathbf{z}_d, \beta_{1:K})] = \\ & = \sum_{i=1}^K \sum_{j=1}^{|\mathcal{V}|} n_{dj} \phi_{dji} \left[\log \Psi(\lambda_{ij}) - \Psi \left(\sum_{s=1}^{|\mathcal{V}|} \lambda_{is} \right) \right]. \end{aligned} \quad (116)$$

5. $E_q[\log q(\beta_{1:K} | \lambda)]$

Con passaggi identici a quelli già visti per il valore atteso del punto 1, cambiando solo ciò che va cambiato:

$$\begin{aligned} & E_q [\log q(\beta_{1:K} | \eta)] = \\ & = \log \Gamma \left(\sum_{j=1}^{|\mathcal{V}|} \lambda_{ij} \right) - \sum_{j=1}^{|\mathcal{V}|} \log \Gamma(\lambda_{ij}) + \\ & + \sum_{i=1}^K \sum_{j=1}^{|\mathcal{V}|} (\lambda_{ij} - 1) \left[\Psi(\lambda_{ij}) - \Psi \left(\sum_{s=1}^{|\mathcal{V}|} \lambda_{is} \right) \right]. \end{aligned} \quad (117)$$

6. $E_q \log [q(\theta_d | \gamma_d)]$

Anche in questo dobbiamo cambiare solo ciò che va cambiato:

$$\begin{aligned} & E_q \log [q(\theta_d | \alpha)] = \\ & = \log \Gamma \left(\sum_{j=1}^K \gamma_{dj} \right) - \sum_{i=1}^K \log \Gamma(\gamma_{di}) + \\ & + \sum_{i=1}^K (\gamma_{di} - 1) \left[\Psi(\gamma_{di}) - \Psi \left(\sum_{j=1}^K \gamma_{dj} \right) \right]. \end{aligned} \quad (118)$$

7. $E_q \log [q(\mathbf{z}_d | \phi_d)]$

Infine, anche per quest'ultimo termine, poche modifiche rispetto a quanto già visto:

$$E_q \log [q(\mathbf{z}_d | \phi_d)] = \sum_{i=1}^K \sum_{j=1}^{|V|} n_{dj} \phi_{dji} \log \phi_{dji}. \quad (119)$$

Raccogliamo tutti i termini dell'energia libera che abbiamo esplicitato:

$$\begin{aligned} \mathcal{L}(\gamma, \phi, \lambda) &= \log \Gamma(\eta | V|) - |V| \log \Gamma(\eta) - \\ &- (\eta - 1) \sum_{j=1}^K \sum_{i=1}^{|V|} \left[\Psi(\lambda_{ij}) - \Psi \left(\sum_{s=1}^{|V|} \lambda_{is} \right) \right] + \\ &+ \sum_{d=1}^M \{ \log \Gamma(K\alpha) - K \log \Gamma(\alpha) \\ &+ \sum_{i=1}^K (\alpha - 1) \left[\Psi(\gamma_{di}) - \Psi \left(\sum_{j=1}^K \gamma_{dj} \right) \right] + \\ &+ \sum_{i=1}^K \sum_{j=1}^{|V|} n_{dj} \phi_{dji} \left[\Psi(\gamma_{di}) - \Psi \left(\sum_{j=1}^K \gamma_{dj} \right) \right] + \\ &+ \sum_{i=1}^K \sum_{j=1}^{|V|} n_{dj} \phi_{dji} \left[\log \Psi(\lambda_{ij}) - \Psi \left(\sum_{s=1}^{|V|} \lambda_{is} \right) \right] \} - \\ &- \log \Gamma \left(\sum_{j=1}^{|V|} \lambda_{ij} \right) + \sum_{j=1}^{|V|} \log \Gamma(\lambda_{ij}) \\ &- \sum_{i=1}^K \sum_{j=1}^{|V|} (\lambda_{ij} - 1) \left[\Psi(\lambda_{ij}) - \Psi \left(\sum_{s=1}^{|V|} \lambda_{is} \right) \right] - \\ &- \sum_{d=1}^M \left\{ \log \Gamma \left(\sum_{j=1}^K \gamma_{dj} \right) - \sum_{i=1}^K \log \Gamma(\gamma_{di}) + \right. \\ &+ \sum_{i=1}^K (\gamma_{di} - 1) \left[\Psi(\gamma_{di}) - \Psi \left(\sum_{j=1}^K \gamma_{dj} \right) \right] - \\ &\left. - \sum_{i=1}^K \sum_{j=1}^{|V|} n_{dj} \phi_{dji} \log \phi_{dji} \right\}. \quad (120) \end{aligned}$$

Poiché gli iperparametri α ed η sono stati fissati, nell'algoritmo EM variazionale per la stima dei parametri variazionali conviene distinguere tra parametri locali, ossia $\{\gamma_d, \phi_d\}$, e il parametro globale λ . Abbiamo quindi i due step:

- **E-step.** Assegnato λ^* massimizziamo l'energia libera documento per documento rispetto ai parametri variazionali locali, ossia:

$$\{\gamma_d^*, \phi_d^*\} = \arg \max_{\{\gamma_d, \phi_d\}} \mathcal{L}(\gamma, \phi, \lambda).$$

- **M-step.** Assegnati $\{\gamma_d^*, \phi_d^*\}$ massimizziamo l'energia libera rispetto a λ :

$$\lambda^* = \operatorname{argmax}_{\lambda} \mathcal{L}(\gamma^*, \phi^*, \lambda).$$

I calcoli per ottenere i punti di massimo locale sono praticamente simili a quelli che abbiamo già visto nella Sezione 5. Ad esempio, tenendo conto degli ovvi vincoli sui parametri variazionali Multinomiali, la relativa Lagrangiana ha espressione:

$$\begin{aligned} \mathcal{L}_{[\phi_{dji}]} &= n_{dj} \phi_{dji} \left[\Psi(\gamma_i) - \Psi \left(\sum_{j=1}^K \gamma_j \right) \right] + \\ &+ n_{dji} \phi_{dji} \left[\Psi(\lambda_{ij}) - \Psi \left(\sum_{s=1}^{|\mathcal{V}|} \lambda_{is} \right) \right] + \\ &+ n_{dji} \phi_{dji} \log \phi_{dji} - \zeta \left[\sum_{i=1}^K \phi_{dji} - 1 \right], \end{aligned} \quad (121)$$

dove ζ indica il moltiplicatore di Lagrange. Prendendo le derivate parziali prime ed uguagliandole a zero otteniamo:

$$\begin{aligned} \frac{\partial \mathcal{L}_{[\phi_{dji}]} }{\partial \phi_{dji}} &= n_{dj} \left[\Psi(\gamma_i) - \Psi \left(\sum_{j=1}^K \gamma_j \right) \right] + \\ &+ n_{dji} \left[\Psi(\lambda_{ij}) - \Psi \left(\sum_{s=1}^{|\mathcal{V}|} \lambda_{is} \right) \right] - \\ &- n_{dj} \log \phi_{dji} - n_{dj} \frac{\phi_{dij}}{\phi_{dji}} - \zeta = 0, \end{aligned} \quad (122)$$

e quindi:

$$\begin{aligned} \phi_{dji} &= \exp(-1 - \zeta) \exp\left(\left[\Psi(\gamma_i) - \Psi\left(\sum_{j=1}^K \gamma_j\right)\right]\right) \times \\ &\times \exp\left(\left[\Psi(\lambda_{ij}) - \Psi\left(\sum_{s=1}^{|\mathcal{V}|} \lambda_{is}\right)\right]\right). \end{aligned} \quad (123)$$

Naturalmente, sappiamo che:

$$\exp\left(\left[\Psi(\gamma_i) - \Psi\left(\sum_{j=1}^K \gamma_j\right)\right]\right) = E_q[\log \theta_{di}], \quad (124)$$

e, allo stesso modo:

$$\exp\left(\left[\Psi(\lambda_{ij}) - \Psi\left(\sum_{s=1}^{|\mathcal{V}|} \lambda_{is}\right)\right]\right) = E_q[\log \beta_{ij}], \quad (125)$$

e quindi possiamo scrivere:

$$\phi_{dji} \propto \exp(E_q[\log \theta_{di}] + E_q[\log \beta_{ij}]), \quad (126)$$

e, naturalmente, l'insieme dei parametri $\{\phi_{dji}\}$ deve essere normalizzato ad 1 sommando rispetto all'indicatore del topic i .

Per quanto riguarda l'altro parametro variazionale locale, ossia γ_d , basta seguire dei passaggi praticamente identici a quelli già visti nel caso precedente per arrivare alla conclusione che l'equazione di aggiornamento (documento per documento) è data da:

$$\gamma_{di} = \alpha + \sum_{j=1}^{|\mathcal{V}|} n_{dj} \phi_{dji}. \quad (127)$$

Altrettanto semplice è verificare che l'equazione di aggiornamento per l'altro parametro variazionale:

$$\lambda_{ij} = \eta + \sum_{d=1}^M n_{dj} \phi_{dji}. \quad (128)$$

Si noti l'evidente simmetria con la (127). Globalmente, l'algoritmo variazionale può essere formulato nel modo seguente (fissata una tolleranza per il monitoraggio dell'energia libera, ad esempio $\epsilon = 10^{-4}$):

Figura 8. Pseudocodice dell'algoritmo EM variazionale per la versione full Bayes del modello LDA.

```

1  inizializza gli elementi di  $\lambda$  in modo casuale
2  while incremento relativo di  $\mathcal{L}(\gamma, \phi, \lambda) > \epsilon$ 
3  E-step
4    for  $d=1$  to  $M$  do
5      inizializza  $\gamma_d^{(0)} \leftarrow 1$  (costante arbitraria)
6      while  $\frac{1}{K} \sum_{i=1}^K |\gamma_{di}| < \epsilon$ 
7        for  $j=1$  to  $|V|$  do
8          for  $i=1$  to  $K$  do
9             $\phi_{dji}^{(t+1)} \leftarrow \exp(\mathbb{E}_q[\log \theta_{di}^{(t)}] + \mathbb{E}_q[\log \beta_{ij}^{(t)}])$ 
10            $\gamma_{di}^{(t+1)} = \alpha + \sum_{j=1}^{|V|} n_{dj} \phi_{dji}^{(t)}$ 
11          end for
12          normalizza la somma  $\phi_{dji}^{(t+1)}$  ad 1 (rispetto all'indicatore del topic  $i$ )
13        end for
14      end for
15 M-step
16    for  $j=1$  to  $|V|$  do
17      for  $i=1$  to  $K$  do
18         $\lambda_{ij}^{(t+1)} = \eta + \sum_{d=1}^M n_{dj} \phi_{dji}^{(t+1)}$ 
19      end for
20    end for
21  until convergence

```

Le stime finali dei parametri differiscono rispetto al caso precedente solo per β_{ij} , poiché per questo modello full Bayes sfruttando il fatto che la distribuzione variazionale su ogni componente di $\beta_{1:K}$ è una Dirichlet:

$$\hat{\beta}_{ij} = \frac{\lambda_{ij}^*}{\sum_{s=1}^{|V|} \lambda_{is}^*}. \quad (129)$$

7. Gibbs sampling collassato per il modello full Bayes

Per il modello LDA full Bayes è possibile utilizzare un campionamento di Gibbs particolarmente efficiente, che campiona esclusivamente dalla distribuzione a posteriori marginale degli indicatori latenti dei topic. La distribuzione **full conditional** dell'indicatore latente z_{dn} per il termine in posizione n del documento d è la seguente:

$$p(z_{dn} | \mathbf{z}_{-(dn)}, \alpha, \eta, \mathcal{D}) = \frac{p(z_{dn}, \mathbf{z}_{-(dn)}, \mathcal{D} | \alpha, \eta)}{p(\mathbf{z}_{-(dn)}, \mathcal{D} | \alpha, \eta)}, \quad (130)$$

dove l'indice $-(dn)$ ha l'ovvio significato che stiamo considerando **tutti** gli altri indicatori latenti dei topic, escluso quello per il documento d in posizione n . Focalizziamoci ora sul numeratore della full conditional, osservando che possiamo scrivere:

$$\begin{aligned} p(z_{dn}, \mathbf{z}_{-(dn)}, \mathcal{D} | \alpha, \eta) &= p(\mathbf{z}_{\mathcal{D}}, \mathcal{D} | \alpha, \eta) = \\ &= \int p(\mathbf{z}_{\mathcal{D}}, \mathcal{D}, \theta_{\mathcal{D}}, \beta_{1:K} | \alpha, \eta) d\theta_{\mathcal{D}} d\beta_{1:K}, \end{aligned} \quad (131)$$

dove $\mathbf{z}_{\mathcal{D}} = \{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_M\}$, con identico significato per $\theta_{\mathcal{D}}$. Si noti che nell'ultimo integrale la full conditional relativa alla sola variabile z_{dn} è ottenuta marginalizzando rispetto ai parametri multinomiali $\theta_{\mathcal{D}}$ e $\beta_{1:K}$, che pertanto non rientrano più in alcun modo nelle iterazioni del campionamento di Gibbs. Questo particolare variante del campionamento di Gibbs, nel quale un insieme di parametri del modello sono stati marginalizzati (e quindi eliminati) dalla full conditional relativa alla variabile di interesse, prende il nome di **Gibbs sampling collassato** (collapsed Gibbs sampling; si veda Liu, 1994, per gli aspetti teorici, nonché Griffiths e Steyvers, 2004).

Sfruttando le relazioni di indipendenza condizionale esistenti tra gli elementi del modello otteniamo:

$$\begin{aligned} &\int p(\mathbf{z}_{\mathcal{D}}, \mathcal{D}, \theta_{\mathcal{D}}, \beta_{1:K} | \alpha, \eta) d\theta_{\mathcal{D}} d\beta_{1:K} = \\ &= \int p(\beta_{1:K} | \eta) p(\theta_{\mathcal{D}} | \alpha) p(\mathbf{z}_{\mathcal{D}} | \theta_{\mathcal{D}}) p(\mathcal{D} | \mathbf{z}_{\mathcal{D}}, \beta_{1:K}) d\theta_{\mathcal{D}} d\beta_{1:K} = \\ &= \int p(\theta_{\mathcal{D}} | \alpha) p(\mathbf{z}_{\mathcal{D}} | \theta_{\mathcal{D}}) d\theta_{\mathcal{D}} \times \\ &\times \int p(\mathcal{D} | \mathbf{z}_{\mathcal{D}}, \beta_{1:K}) p(\beta_{1:K} | \eta) d\beta_{1:K}. \end{aligned} \quad (132)$$

Procediamo sviluppando separatamente i due fattori che compongono l'integrale nella (132). Preliminarmente, osserviamo che per avere il massimo livello di generalità supporremo che gli iperparametri prefissati α ed η non siano costanti, bensì che $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_K)$, mentre $\eta = \{\eta_{ij}\} \in \mathbb{R}^{K \times |V|}$. Naturalmente, nelle applicazioni reali non abbiamo quasi mai informazione sufficiente per settare questi iperparametri ad un tale livello di dettaglio: per tale motivo, al termine della trattazione, forniremo delle espressioni opportunamente semplificate dei risultati finali, nelle quali α ed η ritornano a essere delle costanti, così come nella Sezione precedente. Per il primo fattore della (132) otteniamo:

$$\begin{aligned} & \int p(\theta_{\mathcal{D}}|\alpha)p(\mathbf{z}_{\mathcal{D}}|\theta_{\mathcal{D}})d\theta_{\mathcal{D}} = \\ &= \int \prod_{d=1}^M p(\theta_d|\alpha)p(\mathbf{z}_d|\theta_d)d\theta_1 d\theta_2 \dots d\theta_M = \\ &= \prod_{d=1}^M \int p(\theta_d|\alpha)p(\mathbf{z}_d|\theta_d)d\theta_d, \end{aligned} \quad (133)$$

dove l'ultima uguaglianza è una conseguenza immediata del teorema di Fubini. Ma sappiamo che:

$$p(\theta_d|\alpha) = \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{i=1}^K \theta_{di}^{\alpha_i-1}, \quad (134)$$

mentre:

$$p(\mathbf{z}_d) = \prod_{n=1}^{N_d} p(z_{dn}|\theta_d) = \prod_{n=1}^{N_d} \prod_{i=1}^K \theta_{di}^{z_{dn}^i}. \quad (135)$$

Riscriviamo quest'ultima nel modo seguente:

$$\prod_{n=1}^{N_d} \prod_{i=1}^K \theta_{di}^{z_{dn}^i} = \prod_{i=1}^K \prod_{n=1}^{N_d} \theta_{di}^{z_{dn}^i} = \prod_{i=1}^K \theta_{di}^{\sum_{n=1}^{N_d} z_{dn}^i}, \quad (136)$$

ponendo:

$$C(d, i) \equiv \sum_{n=1}^{n_d} z_n^i, \quad (137)$$

dove è ovvio che $C(d, i)$ indica il **numero di termini nel documento d -esimo assegnati al topic i -esimo**. Con questa posizione, otteniamo:

$$\begin{aligned} & \prod_{d=1}^M \int p(\theta_d | \alpha) p(\mathbf{z}_d | \theta_d) d\theta_d = \\ &= \prod_{d=1}^M \int \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{i=1}^K \theta_{di}^{\alpha_i - 1} \prod_{i=1}^K \theta_{di}^{C(d,i)} d\theta_d = \\ &= \prod_{d=1}^M \int \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{i=1}^K \theta_{di}^{\alpha_i + C(d,i) - 1} d\theta_d. \end{aligned} \quad (138)$$

Lo sviluppo del secondo termine dell'integrale (132) è abbastanza simile:

$$\begin{aligned} & \int p(\mathcal{D} | \mathbf{z}_D, \beta_{1:K}) p(\beta_{1:K} | \eta) d\beta_{1:K} = \\ &= \prod_{i=1}^K \int p(\beta_i | \eta) \prod_{d=1}^M \prod_{n=1}^{N_d} p(w_{dn} | \beta_{1:K}, z_{dn}) d\beta_i. \end{aligned} \quad (139)$$

Ma:

$$p(\beta_i | \eta) = \frac{\Gamma(\sum_{j=1}^{|\mathcal{V}|} \eta_{ij})}{\prod_{j=1}^{|\mathcal{V}|} \Gamma(\eta_{ij})} \prod_{j=1}^{|\mathcal{V}|} \beta_{ij}^{\eta_{ij} - 1}, \quad (140)$$

mentre per il secondo fattore, considerando che esso si trova all'interno di una produttoria sull'indicatore i dei topic (ossia per l' i -esimo termine di quest'ultima produttoria l'indicatore latente z_{dn} assegna il termine in posizione n nel documento d all' i -esimo topic), possiamo evidentemente scrivere:

$$\begin{aligned}
\prod_{d=1}^M \prod_{n=1}^{N_d} p(w_{dn} | \beta_{1:K}, z_{dn}) &= \prod_{d=1}^M \prod_{n=1}^{N_d} \prod_{j=1}^{|V|} \beta_{ij}^{w_n^j} = \\
&= \prod_{j=1}^{|V|} \prod_{d=1}^M \prod_{n=1}^{N_d} \beta_{ij}^{w_n^j} = \prod_{j=1}^{|V|} \beta_{ij}^{\sum_{d=1}^M \sum_{n=1}^{N_d} w_n^j}.
\end{aligned} \tag{141}$$

e anche in questo caso poniamo:

$$C(i, j) = \sum_{d=1}^M \sum_{n=1}^{N_d} w_n^j, \tag{142}$$

dove $C(i, j)$ è evidentemente uguale al **numero complessivo di termini assegnati al topic i -esimo che coincidono con il j -esimo token del vocabolario V** . Abbiamo, pertanto:

$$\begin{aligned}
&\prod_{i=1}^K \int p(\beta_i | \eta) \prod_{d=1}^M \prod_{n=1}^{N_d} p(w_{dn} | \beta_{1:K}, z_{dn}) d\beta_i = \\
&= \prod_{i=1}^K \int \frac{\Gamma(\sum_{j=1}^{|V|} \eta_{ij})}{\prod_{j=1}^{|V|} \Gamma(\eta_{ij})} \prod_{j=1}^{|V|} \beta_{ij}^{\eta_{ij}-1} \prod_{j=1}^{|V|} \beta_{ij}^{C(i,j)} d\beta_i = \\
&= \prod_{i=1}^K \int \frac{\Gamma(\sum_{j=1}^{|V|} \eta_{ij})}{\prod_{j=1}^{|V|} \Gamma(\eta_{ij})} \prod_{j=1}^{|V|} \beta_{ij}^{\eta_{ij}+C(i,j)-1} d\beta_i.
\end{aligned} \tag{143}$$

Possiamo ora rimettere insieme le espressioni dei due fattori che abbiamo esplicitato, e fare tutte le dovute semplificazioni:

$$\begin{aligned}
& \prod_{d=1}^M \int \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{i=1}^K \theta_{di}^{\alpha_i + C(d,i) - 1} d\theta_d \times \\
& \times \prod_{i=1}^K \int \frac{\Gamma(\sum_{j=1}^{|V|} \eta_{ij})}{\prod_{j=1}^{|V|} \Gamma(\eta_{ij})} \prod_{j=1}^{|V|} \beta_{ij}^{\eta_{ij} + C(i,j) - 1} d\beta_i = \\
& = \prod_{d=1}^M \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \frac{\prod_{i=1}^K \Gamma(C(d,i) + \alpha_i)}{\Gamma(\sum_{i=1}^K (C(d,i) + \alpha_i))} \times \\
& \quad \underbrace{\int \frac{\Gamma(\sum_{i=1}^K (C(d,i) + \alpha_i))}{\prod_{i=1}^K \Gamma(C(d,i) + \alpha_i)} \prod_{i=1}^K \theta_{di}^{\alpha_i + C(d,i) - 1} d\theta_d}_{\equiv 1} \times \\
& \quad \times \prod_{i=1}^K \frac{\Gamma(\sum_{j=1}^{|V|} \eta_{ij})}{\prod_{j=1}^{|V|} \Gamma(\eta_{ij})} \frac{\prod_{j=1}^{|V|} \Gamma(C(i,j) + \eta_{ij})}{\Gamma(\sum_{j=1}^{|V|} (C(i,j) + \eta_{ij}))} \times \\
& \quad \underbrace{\int \frac{\Gamma(\sum_{j=1}^{|V|} (C(i,j) + \eta_{ij}))}{\prod_{j=1}^{|V|} \Gamma(C(i,j) + \eta_{ij})} \prod_{j=1}^{|V|} \beta_{ij}^{\eta_{ij} + C(i,j) - 1} d\beta_i}_{\equiv 1} = \\
& = \prod_{d=1}^M \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \frac{\prod_{i=1}^K \Gamma(C(d,i) + \alpha_i)}{\Gamma(\sum_{i=1}^K (C(d,i) + \alpha_i))} \times \\
& \quad \times \prod_{i=1}^K \frac{\Gamma(\sum_{j=1}^{|V|} \eta_{ij})}{\prod_{j=1}^{|V|} \Gamma(\eta_{ij})} \frac{\prod_{j=1}^{|V|} \Gamma(C(i,j) + \eta_{ij})}{\Gamma(\sum_{j=1}^{|V|} (C(i,j) + \eta_{ij}))} \propto \\
& \propto \text{eliminiamo termini cost. che dipendono solo da } \alpha \text{ ed } \eta \propto \\
& \propto \prod_{d=1}^M \frac{\prod_{i=1}^K \Gamma(C(d,i) + \alpha_i)}{\Gamma(\sum_{i=1}^K (C(d,i) + \alpha_i))} \times \prod_{i=1}^K \frac{\prod_{j=1}^{|V|} \Gamma(C(i,j) + \eta_{ij})}{\Gamma(\sum_{j=1}^{|V|} (C(i,j) + \eta_{ij}))}.
\end{aligned} \tag{144}$$

Possiamo semplificare ulteriormente queste espressioni introducendo una nuova notazione vettorializzata. Scriviamo:

$$\begin{aligned}
C_d & \equiv (C(d, 1), C(d, 2), \dots, C(d, K)) \\
B(\alpha + C_d) & \equiv \frac{\prod_{i=1}^K \Gamma(C(d, i) + \alpha_i)}{\Gamma(\sum_{i=1}^K (C(d, i) + \alpha_i))}.
\end{aligned} \tag{145}$$

Allo stesso modo, se poniamo:

$$C_i = (C(i, 1), C(i, 2), \dots, C(i, |V|)), \quad (146)$$

possiamo definire:

$$B(\eta_i + C_i) \equiv \frac{\prod_{j=1}^{|V|} \Gamma(C(i, j) + \eta_{ij})}{\Gamma\left(\sum_{j=1}^{|V|} (C(i, j) + \eta_{ij})\right)}. \quad (147)$$

Con questa notazione, il prodotto dei due fattori che compaiono nell'ultima riga della (144) può essere semplicemente riscritto come:

$$\left[\prod_{d=1}^M B(\alpha + C_d) \right] \left[\prod_{i=1}^K B(\eta_i + C_i) \right]. \quad (148)$$

Dopo questo lungo excursus sul numeratore della full conditional per z_{dn} , volto ad esplicitarne l'espressione, consideriamo il denominatore, osservando che possiamo scrivere (vale ovviamente la stessa marginalizzazione rispetto ai parametri Multinomiali, già utilizzata per il numeratore):

$$\begin{aligned} & p(\mathbf{z}_{-(dn)}, \mathcal{D} | \alpha, \eta) = \\ & = p(\mathbf{z}_{-(dn)}, w_{dn}, \mathbf{w}_{-(dn)} | \alpha, \eta) = \\ & = p(\mathbf{w}_{-(dn)} | \mathbf{z}_{-(dn)}, \alpha, \eta) p(\mathbf{z}_{-(dn)} | \alpha, \eta) p(w_{dn} | \alpha, \eta) \propto \\ & \propto p(\mathbf{w}_{-(dn)} | \mathbf{z}_{-(dn)}, \alpha, \eta) p(\mathbf{z}_{-(dn)} | \alpha, \eta) = \\ & = p(\mathbf{w}_{-(dn)}, \mathbf{z}_{-(dn)} | \alpha, \eta). \end{aligned} \quad (149)$$

Questa espressione è formalmente identica a quella del numeratore, si riveda la (131), con la sola differenza che per il numeratore stiamo considerando tutti i termini e tutti i relativi indicatori latenti del topic, **tranne** il termine posizione n nel

documento d . Pertanto, con gli stessi passaggi già visti per il numeratore, il denominatore può essere scritto come:

$$p(\mathbf{w}_{-(dn)}, \mathbf{z}_{-(dn)} | \alpha, \eta) \propto \left[\prod_{d=1}^M B(\alpha + C_d)_{-(dn)} \right] \left[\prod_{i=1}^K B(\eta_i + C_i)_{-(dn)} \right], \quad (150)$$

dove:

$$B(\alpha + C_d)_{-(dn)} \equiv \frac{\prod_{i=1}^K \Gamma(C(d, i)_{-(dn)} + \alpha_i)}{\Gamma\left(\sum_{i=1}^K (C(d, i)_{-(dn)} + \alpha_i)\right)}, \quad (151)$$

con $c(d, i)_{-(dn)}$ che indica il numero di termini nel documento d assegnati all' i -esimo topic, senza tener conto del termine in posizione n . In modo del tutto analogo per l'altro termine:

$$B(\eta_i + C_i)_{-(dn)} = \frac{\prod_{j=1}^{|V|} \Gamma(C(i, j)_{-(dn)} + \eta_{ij})}{\Gamma\left(\sum_{j=1}^{|V|} (C(i, j)_{-(dn)} + \eta_{ij})\right)}, \quad (152)$$

mella quale $C(i, j)_{-(dn)}$, così come nel caso precedente, indica il numero di termini assegnati al topic i -esimo che coincidono con il token j -esimo del vocabolario dei termini, senza tener conto del termine in posizione n del documento d .

Per queste variabili di conteggio nel calcolo delle quali è stato escluso il termine in posizione n del documento d valgono delle ovvie uguaglianze. In particolare:

$$C(d, i) = C(d, i)_{-(dn)}, \quad (153)$$

se il termine in posizione n del documento d non è assegnato al topic i , mentre:

$$C(d, i) = C(d, i)_{-(dn)} + 1, \quad (154)$$

se il termine in posizione n del documento d è assegnato al topic i . Simili uguaglianze sono valide anche per l'altra variabile di conteggio, ossia:

$$C(i, j) = C(i, j)_{-(dn)}, \quad (155)$$

se tra i termini assegnati al topic i , quello in posizione (dn) non coincide con il token j , mentre:

$$C(i, j) = C(i, j) + 1, \quad (156)$$

se tra i termini assegnati al topic i , quello in posizione (dn) coincide con il token j .

Con la notazione vettoriale che abbiamo introdotto, la full conditional (130) può essere riscritta come:

$$\begin{aligned} & p(z_{dn} | \mathbf{z}_{-(dn)}, \alpha, \eta, \mathcal{D}) = \\ &= \frac{p(z_{\mathcal{D}}, \mathcal{D} | \alpha, \eta)}{p(\mathbf{z}_{-(dn)}, \mathcal{D} | \alpha, \eta)} = \\ &= \frac{p(z_{\mathcal{D}}, \mathcal{D} | \alpha, \eta)}{p(\mathbf{w}_{-(dn)}, \mathbf{z}_{-(dn)} | \alpha, \eta)} \propto \\ &\propto \left[\frac{\prod_{d=1}^M B(\alpha + C_d)}{\prod_{d=1}^M B(\alpha + C_d)_{-(dn)}} \right] \left[\frac{\prod_{i=1}^K B(\eta_i + C_i)}{\prod_{i=1}^K B(\eta_i + C_i)_{-(dn)}} \right]. \end{aligned} \quad (157)$$

Consideriamo ora la full conditional per la posizione corrente $(d'n')$ per il topic i' , ossia:

$$p(z_{d'n'} = i' | \mathbf{z}_{-(d'n')}, \alpha, \eta, \mathcal{D}). \quad (158)$$

Naturalmente, l'indicatore del topic è un vettore K -dimensionale, ma in questo caso a scopo di semplificazione abuseremo della notazione supponendo semplicemente che $i' \in \{1, 2, \dots, K\}$. Tutti fattori moltiplicativi che nella full conditional corrente non dipendono da i' possono essere fatti scomparire nella costante di normalizzazione, ossia:

$$\begin{aligned}
& \frac{\prod_{d=1}^M B(\alpha + C_d)}{\prod_{d=1}^M B(\alpha + C_d)_{-(dmr)}} \propto \\
& \propto \frac{B(\alpha + C_{d'})}{B(\alpha + C_{d'})_{-(dmr)}} = \\
& = \frac{\prod_{i=1}^K \Gamma(C(d', i) + \alpha_i)}{\Gamma(\sum_{i=1}^K (C(d', i) + \alpha_i))} \times \\
& \times \left[\frac{\prod_{i=1}^K \Gamma(C(d', i)_{-(dmr)} + \alpha_i)}{\Gamma(\sum_{i=1}^K (C(d', i)_{-(dmr)} + \alpha_i))} \right]^{-1} = \\
& = \frac{\prod_{i=1}^K \Gamma(C(d', i) + \alpha_i)}{\prod_{i=1}^K \Gamma(C(d', i)_{-(dmr)} + \alpha_i)} \times \\
& \times \frac{\Gamma(\sum_{i=1}^K (C(d', i)_{-(dmr)} + \alpha_i))}{\Gamma(\sum_{i=1}^K (C(d', i) + \alpha_i))} \propto \\
& \propto \text{elim. i termini che non dipend. da } i' \propto \\
& \propto \frac{\Gamma(C(d', i') + \alpha_{i'})}{\Gamma(C(d', i')_{-(dmr)} + \alpha_{i'})} \times \\
& \times \frac{\Gamma(\sum_{i=1}^K (C(d', i)_{-(dmr)} + \alpha_i))}{\Gamma(\sum_{i=1}^K (C(d', i) + \alpha_i))} =
\end{aligned} \tag{159}$$

Osserviamo ora che, grazie alla (153) e alla (154):

$$\begin{aligned}
& \frac{\Gamma(C(d', i') + \alpha_{i'})}{\Gamma(C(d', i')_{-(dmr)} + \alpha_{i'})} = \\
& = \frac{\Gamma(C(d', i')_{-(dmr)} + \alpha_{i'} + 1)}{\Gamma(C(d', i')_{-(dmr)} + \alpha_{i'})},
\end{aligned} \tag{160}$$

e allo stesso modo:

$$\begin{aligned}
& \frac{\Gamma(\sum_{i=1}^K (C(d', i)_{-(dmr)} + \alpha_i))}{\Gamma(\sum_{i=1}^K (C(d', i) + \alpha_i))} = \\
& = \frac{\Gamma(\sum_{i \neq i'} (C(d', i)_{-(dmr)} + \alpha_j) + [C(d', i')_{-(dmr)} + \alpha_{i'}])}{\Gamma(\sum_{i \neq i'} (C(d', i)_{-(dmr)} + \alpha_i) + [C(d', i')_{-(dmr)} + \alpha_{i'}] + 1)}.
\end{aligned} \tag{161}$$

Se utilizziamo la ben nota relazione $\Gamma(x + 1) = x\Gamma(x)$, il risultato ottenuto nella (160) diventa:

$$\begin{aligned} & \frac{\Gamma(C(d', i')_{-(dmv)} + \alpha_{i'} + 1)}{\Gamma(C(d', i')_{-(dmv)} + \alpha_{i'})} = \\ & = \frac{(C(d', i')_{-(dmv)} + \alpha_{i'}) \times \Gamma(C(d', i')_{-(dmv)} + \alpha_{i'})}{\Gamma(C(d', i')_{-(dmv)} + \alpha_{i'})} = \\ & = (C(d', i')_{-(dmv)} + \alpha_{i'}), \end{aligned} \quad (162)$$

e sfruttando la stessa relazione, è immediato che anche la (161) si semplifica in:

$$\begin{aligned} & \frac{\Gamma(\sum_{i \neq i'} (C(d', i)_{-(dmv)} + \alpha_i) + [C(d', i')_{-(dmv)} + \alpha_{i'}])}{\Gamma(\sum_{i \neq i'} (C(d', i)_{-(dmv)} + \alpha_i) + [C(d', i')_{-(dmv)} + \alpha_{i'}] + 1)} = \\ & = \frac{1}{\sum_{i=1}^K [C(d', i)_{-(dmv)} + \alpha_i]}. \end{aligned} \quad (163)$$

Per semplificare il secondo fattore nella (157) sono validi passaggi assolutamente analoghi; naturalmente, la posizione corrente sarà $(d'n')$ e il token corrente j' , e quindi il fattore in questione si riduce a:

$$\frac{(C(i', j')_{-(dmv)} + \eta_{i'j'})}{\sum_{j=1}^{|V|} [C(i', j)_{-(dmv)} + \eta_{i'j}]} \quad (164)$$

Siamo quasi arrivati all'espressione finale, in quanto grazie a ciò che abbiamo dimostrato possiamo scrivere:

$$\begin{aligned} & p(z_{d'n'} = i' | \mathbf{z}_{-(d'n')}, \alpha, \eta, \mathcal{D}) \propto \\ & \propto \left[\frac{(C(d', i')_{-(dmv)} + \alpha_{i'})}{\sum_{i=1}^K [C(d', i)_{-(dmv)} + \alpha_i]} \right] \left[\frac{(C(i', j')_{-(dmv)} + \eta_{i'j'})}{\sum_{j=1}^{|V|} [C(i', j)_{-(dmv)} + \eta_{i'j}]} \right]. \end{aligned} \quad (165)$$

Resta solo da osservare che se la posizione corrente è $(d'n')$ la quantità:

$$\sum_{i=1}^K [C(d', i)_{-(dmv)} + \alpha_i], \quad (166)$$

è costante al variare dell'indicatore del topic i' , e quindi possiamo scrivere:

$$\begin{aligned}
& p(z_{dn} = i' | \mathbf{z}_{-(dn)}, \alpha, \eta, \mathcal{D}) \propto \\
& \propto \frac{[(C(d', i')_{-(dn)} + \alpha_{i'})][C(i', j')_{-(dn)} + \eta_{i'j'}]}{\sum_{j=1}^{|V|} [C(i', j)_{-(dn)} + \eta_{ij}]} .
\end{aligned} \tag{167}$$

Se gli iperparametri α ed η sono costanti (distribuzioni di Dirichlet simmetriche), questa espressione si riduce ulteriormente a:

$$\begin{aligned}
& p(z_{dn} = i' | \mathbf{z}_{-(dn)}, \alpha, \eta, \mathcal{D}) \propto \\
& \propto \frac{[(C(d', i')_{-(dn)} + \alpha)][C(i', j')_{-(dn)} + \eta]}{\sum_{j=1}^{|V|} (C(i', j)_{-(dn)} + \eta|V|)} .
\end{aligned} \tag{168}$$

In entrambi i casi, queste probabilità non normalizzate andranno normalizzate ad 1, dividendo per la relativa somma rispetto all'indicatore del topic i' .

Ad ogni iterazione dell'algoritmo di Gibbs collassato che abbiamo costruito, le assegnazioni dei termini ai topic vengono aggiornate. All'istante iniziale, per ciascun termine w_{dn} gli indicatori dei topic vengono inizializzati in modo causale con un valore compreso tra 1 e K . Dopo un numero di iterazioni sufficiente, la catena raggiungerà come distribuzione target la distribuzione marginale a posteriori degli indicatori dei topic latenti. Ciascuna full conditional tende a concentrarsi su un sottoinsieme limitato di possibili topic, e le assegnazioni dei topic ai singoli termini tendono a diventare stabili.

Una volta che abbiamo a disposizione una campione $\mathbf{z}_{\mathcal{D}}^{(\text{Gibbs})}$ dalla distribuzione marginale a posteriori di $\mathbf{z}_{\mathcal{D}}$ come possiamo utilizzarlo? Per prima cosa, calcoliamo il valore corrente delle quantità $C(d, i)$ e $C(i, j)$ per ogni possibile posizione (dn) . Indichiamo tali conteggi rispettivamente con $C^{(\text{Gibbs})}(d, i)$ e $C^{(\text{Gibbs})}(i, j)$.

Se gli indicatori latenti dei topic sono noti, l'inferenza a posteriori sui parametri Multinomiali può essere effettuata utilizzando anziché la distribuzione a posteriori marginale di ciascun parametro (intrattabile), la distribuzione a posteriori marginale di ciascun parametro condizionale agli indicatori dei topic (noti). Come vedremo tra breve, questa distribuzione è facilmente trattabile grazie alla coniugazione Dirichlet-Multinomiale. Pertanto, se consideriamo il parametro vettoriale θ_d il nostro obiettivo è la distribuzione:

$$p(\theta_{\mathcal{D}} | \mathbf{z}_{\mathcal{D}}, \mathcal{D}), \tag{169}$$

che può essere scritta come (a questo punto sopprimiamo gli iperparametri costanti):

$$\begin{aligned}
 p(\theta_{\mathcal{D}}|\mathbf{z}_{\mathcal{D}}, \mathcal{D}) &\propto \\
 &\propto p(\theta_{\mathcal{D}})p(\mathbf{z}_{\mathcal{D}}|\theta_{\mathcal{D}})p(\mathcal{D}|\beta_{1:K}, \mathbf{z}_{\mathcal{D}}) \propto \\
 &\propto p(\theta_{\mathcal{D}})p(\mathbf{z}_{\mathcal{D}}|\theta_{\mathcal{D}}).
 \end{aligned} \tag{170}$$

Pertanto, documento per documento:

$$\begin{aligned}
 p(\theta_d)p(\mathbf{z}_d|\theta_d) &= p(\theta_d) \prod_{n=1}^{N_d} p(z_{dn}|\theta_d) \propto \\
 &\propto \prod_{i=1}^K \theta_{di}^{\alpha-1} \prod_{n=1}^{N_d} \prod_{i=1}^K \theta_{di}^{z_{dn}^i} = \\
 &= \prod_{i=1}^K \theta_{di}^{\alpha-1} \prod_{i=1}^K \theta_{di}^{\sum_{n=1}^{N_d} z_{dn}^i} = \\
 &= \prod_{i=1}^K \theta_{di}^{\alpha-1} \prod_{i=1}^K \theta_{di}^{C(d,i)} = \\
 &= \prod_{i=1}^K \theta_{di}^{\alpha+C(d,i)-1} \propto \\
 &\propto \text{Dirichlet}_K(\alpha + C_d),
 \end{aligned} \tag{171}$$

e, pertanto, la stima Bayesiana a posteriori di θ_{di} ottenuta da questa distribuzione è:

$$\begin{aligned}
 E(\theta_{di}|\mathbf{z}_{\mathcal{D}}, \mathcal{D}) &= \frac{C(d,i) + \alpha}{\sum_{j=1}^K (C(d,j) + \alpha)} = \\
 &= \frac{C(d,i) + \alpha}{\sum_{j=1}^K (C(d,j)) + K\alpha},
 \end{aligned} \tag{172}$$

che può essere ovviamente approssimata come:

$$\begin{aligned}
 \hat{\theta}_{di}^G &= \frac{C(d,i)^{(\text{Gibbs})} + \alpha}{\sum_{j=1}^K (C(d,j)^{(\text{Gibbs})} + \alpha)} = \\
 &= \frac{C(d,i)^{(\text{Gibbs})} + \alpha}{\sum_{j=1}^K (C(d,j))^{(\text{Gibbs})} + K\alpha}.
 \end{aligned} \tag{173}$$

Per β_{ij} vale una procedura assolutamente analoga, osservando che:

$$p(\beta_{1:K} | \mathbf{z}_D, \mathcal{D}) \propto p(\beta_{1:K}, \mathbf{z}_D, \mathcal{D}) = p(\mathcal{D} | \beta_{1:K}, \mathbf{z}_D) p(\beta_{1:K}). \quad (174)$$

Sfruttando ancora una volta la coniugazione Dirichlet-Multinomiale con calcoli identici a quelli già visti arriviamo a:

$$\beta_i | \mathbf{z}_D, \mathcal{D} \sim \text{Dirichlet}_{|V|}(\eta + C_i), \quad (175)$$

e quindi:

$$\begin{aligned} E(\beta_{ij} | \mathbf{z}_D, \mathcal{D}) &= \frac{C(i, j) + \eta}{\sum_{s=1}^{|V|} (C(i, s) + \eta)} = \\ &= \frac{C(i, j) + \eta}{\sum_{s=1}^{|V|} (C(i, s)) + |V|\eta}, \end{aligned} \quad (176)$$

che ancora una volta può essere approssimato con:

$$\begin{aligned} \hat{\beta}_{ij}^G &= \frac{C(i, j)^{(\text{Gibbs})} + \eta}{\sum_{s=1}^{|V|} (C(i, s)^{(\text{Gibbs})} + \eta)} = \\ &= \frac{C(i, j)^{(\text{Gibbs})} + \eta}{\sum_{s=1}^{|V|} (C(i, s)^{(\text{Gibbs})}) + |V|\eta}. \end{aligned} \quad (177)$$

Una spiegazione equivalente di questa procedura è basata su una tecnica, utilizzata per ridurre la varianza delle stime a posteriori MCMC, nota come Rao-Blackwellizzazione (a riguardo si veda Robert et al., 2018).

Ovviamente, sembrerebbe altrettanto sensato calcolare la (172) e la (176) come medie empiriche (Heinrich, 2008): abbiamo pertanto bisogno, una volta che la convergenza sia stata raggiunta, di conservare L valori corrispondenti ad L stati $\mathbf{z}_D^{(t)}$ campionati dalla catena di Markov (opportunosamente distanziati di un certo intervallo di thinning, al fine di ridurre o l'eliminare del tutto l'autocorrelazione esistente nei valori generati) delle quantità $C^{(t)}(d, i)$ e $C^{(t)}(i, j)$, nonché di:

$$\begin{aligned}
C^{(t)}(d) &= \sum_{j=1}^K C^{(t)}(d, j), \\
C^{(t)}(i) &= \sum_{s=1}^{|V|} C^{(t)}(i, s).
\end{aligned}
\tag{178}$$

Pertanto, per ciascuna delle L iterazioni che abbiamo conservato possiamo calcolare $\hat{\theta}_{di}^{G,(t)}$ e $\hat{\beta}_{ij}^{G,(t)}$ (con $t = 1, 2, \dots, L$), che saranno poi ridotti ad una unica stima mediante una media empirica. Tuttavia, questo modo di procedere è abbastanza rischioso per la classe di modelli che stiamo trattando, poiché la verosimiglianza del modello LDA è invariante per permutazioni arbitrarie dei topic. Ad esempio, la verosimiglianza assume esattamente lo stesso valore se il topic 1, che compare nella prima riga della matrice $\beta_{1:K}$, viene spostato in un'altra riga. Questo sottile problema di non identificabilità è ben noto in letteratura, e prende il nome di **label switching** (Stephens, 2000). Durante la fase di campionamento, l'algoritmo di Gibbs salta continuamente tra queste versioni della verosimiglianza permutazionalmente equivalenti, con la conseguenza che le probabilità $\hat{\beta}_{ij}^{G,(t)}$ che abbiamo stimato potrebbero in realtà riferirsi a distribuzioni (topic) differenti.

Di seguito, abbiamo riportato lo pseudocodice del campionamento di Gibbs collassato: per non appesantire troppo, abbiamo eliminato ogni riferimento all'istante di campionamento t , poiché gli aggiornamenti sono chiari dal contesto. L'implementazione che riportiamo è particolarmente efficiente. Per la posizione corrente (dn) e dato il topic corrente i , tutte le variabili di conteggio vengono decrementate di 1 per poter calcolare la relativa distribuzione full conditional; esse vengono poi nuovamente incrementate di 1 solo se il topic aggiornato coincide con quello precedente. Altrimenti la variabile di conteggio relativa al topic precedente resta decrementata di 1, mentre viene incrementata di 1 solo quella relativa al topic aggiornato i' . L'unica variabile di conteggio che viene sempre prima decrementata di 1 e poi incrementata è $C(d)$: ciò è ovvio, poiché ogni termine deve essere assegnato ad un topic. Anzi, al termine del ciclo for su n (quando tutti gli indicatori dei topic sul documento d sono stati aggiornati), è ovvio che risulterà in ogni caso (qualunque sia la distribuzione corrente delle assegnazioni dei topic tra i termini):

$$\sum_{j=1}^K C^{(t)}(d, j) = C^{(t)}(d) = N_d. \quad (179)$$

Figura 9. Pseudocodice dell'algorithm collapsed Gibbs sampling per la versione full Bayes del modello LDA

```

1  C(d, i), C(i, j), C(i), C(d) ← 0
2  for d ∈ {1, 2, ..., M} do
3    for n ∈ {1, 2, ..., N_d} do
4      inizializza in modo casuale l'indicatore del topic: z_{dn} ~ Multinomial_K(1/K)
5      dato il token j in posizione corrente
6      C(d, i) → C(d, i) + 1
7      C(i, j) → C(i, j) + 1
8      C(i) → C(i) + 1
9      C(d) → C(d) + 1
10   end for
11 end for
12 while la convergenza è raggiunta
13 Gibbs sampling collassato nel periodo transitorio e in quello di burn-in
14 for d ∈ {1, 2, ..., M} do
15   for n ∈ {1, 2, ..., N_d} do
16     dato l'assegnazione corrente al topic i per il termine w_{dn}
17     dato il token corrente j per il termine w_{dn}
18     C(d, i) → C(d, i) - 1
19     C(i, j) → C(i, j) - 1
20     C(i) → C(i) - 1
21     C(d) → C(d) - 1
22     for i ∈ {1, 2, ..., K} do
23       calcola p(z_{dn} = i | z_{-(dn)}, D)
24       normalizza ad 1
25     end for
26     aggiorna l'assegnazione al topic per il termine w_{dn} mediante p(z_{dn} | z_{-(dn)}, D)
27     dato il topic corrente i'
28     C(d, i') → C(d, i') + 1
29     C(i', j) → C(i', j) + 1
30     C(i') → C(i') + 1
31     C(d) → C(d) + 1
32   end for
33 end for
34 until è stato generato un numero sufficiente di stati della catena

```

8. Scelta del numero ottimale di topic

Il modello LDA, nelle sue varie declinazioni, contiene come input fondamentale che deve essere prefissato dall'utente il numero di topic K . Abbiamo pertanto bisogno di

un qualche metodo di valutazione che ci consenta di scegliere, prima di apprendere i parametri del modello, di settare K .

Un modo possibile di procedere è quello di valutare l'accuratezza previsiva per K prefissato. Il problema di scegliere K può, pertanto, essere visto come un problema di valutazione e scelta del modello ottimale dal punto di vista dell'accuratezza previsiva, facendo variare K in un range opportuno: il modello 'migliore' sarà quello corrispondente a quel valore di K che consente di riprodurre nel miglior modo possibile i dati futuri.

In altre parole, da un punto di vista previsivo il 'miglior' modello sarà quello che assegna la probabilità più elevata di verificarsi ad un insieme di documenti che non appartengono al corpus che è stato utilizzato per stimare i parametri del modello stesso. Spesso un insieme di documenti futuri è ottenuto rimuovendo dai dati disponibili una percentuale prefissata di documenti del corpus estratti modo casuale (in modo da non introdurre distorsioni sistematiche):

$$\mathcal{D}^{\text{new}} = \{\mathbf{w}_1^{\text{new}}, \mathbf{w}_2^{\text{new}}, \dots, \mathbf{w}_T^{\text{new}}\}. \quad (180)$$

In questo caso l'insieme \mathcal{D}^{new} è comunemente noto come **test set**, ovvero anche **holdout set**. L'uso di un insieme di documenti distinto da quello utilizzato per stimare i parametri evita di essere eccessivamente ottimisti sulla performance previsiva. Infatti, se gli stessi dati sono utilizzati due volte (tanto per la stima dei parametri quanto per la scelta del modello), verrà semplicemente favorito quel modello che ha migliore capacità di riprodurre i dati correnti utilizzati per l'apprendimento dei parametri, senza però alcuna garanzia di ottimalità quando siano presentati nuovi dati futuri effettivamente mai visti prima.

Per un documento di test, la log-probabilità marginale di osservare il documento d nell'insieme di test ha espressione generica:

$$\log p(\mathbf{w}_d^{\text{new}} | \alpha, \eta). \quad (181)$$

Ovviamente, documenti più lunghi avranno più incertezza: quindi, tanto maggiore è la lunghezza del documento di test corrente, tanto minore è le possibilità del modello corrente di assegnare ad esso una probabilità marginale elevata. Consideriamo, ad esempio, due documenti descritti da $\mathbf{w}_{d_1}^{\text{new}}$ e $\mathbf{w}_{d_2}^{\text{new}}$, per i quali la log-probabilità marginale è pari rispettivamente a $\log 0.01 \approx -4.61$ e $\log 0.0001 \approx -9.21$. Dunque, sotto il modello corrente il documento d_1^{new} ha una probabilità più elevata di verificarsi di d_2^{new} . Tuttavia, se $N_{d_1}^{\text{new}} = 100$ mentre $N_{d_2}^{\text{new}} = 1000$, possiamo normalizzare per la lunghezza del documento (ossia calcolare la log-probabilità per

termine), ottenendo $\log 0.01 \approx -0.0461$ e $\log 0.0001/1000 \approx -0.0009$. Una volta che abbiamo rimosso l'effetto della lunghezza del documento, il vettore $\mathbf{w}_{d_2}^{\text{new}}$ è quello più supportato dal modello corrente.

Pertanto, se abbiamo T documenti nel test set, è naturale aggregare le log-probabilità marginali su scala additiva, e normalizzare rispetto al numero complessivo di termini, ossia:

$$\frac{\sum_{d=1}^T \log p(\mathbf{w}_d^{\text{new}} | \alpha, \eta)}{\sum_{d=1}^T N_d^{\text{new}}}. \quad (182)$$

Infine, una misura di accuratezza comunemente utilizzata è una trasformazione monotona di questa quantità, nota come **perplexità** (perplexity), ossia:

$$\text{perplexity}(\mathcal{D}^{\text{new}}) = \exp \left\{ - \frac{\sum_{d=1}^T \log p(\mathbf{w}_d^{\text{new}} | \alpha, \eta)}{\sum_{d=1}^T N_d^{\text{new}}} \right\}. \quad (183)$$

Evidentemente, questa misura assume valori positivi: inoltre, data la presenza del segno negativo davanti al termine di verosimiglianza, il modello corrente che più supporta dal punto di vista previsivo il corpus di test sarà quello avente la perplexità più piccola.

Il calcolo della (181) è particolarmente complicato, poiché di solito ha a che fare con integrazioni multidimensionali non esplicitabili (questo problema è già stato discusso nella Sezione 2). Osserviamo che la log-verosimiglianza marginale dell'insieme di test \mathcal{D}^{new} può essere scritta come:

$$\begin{aligned} & \log p(\mathcal{D}^{\text{new}} | \alpha, \eta) = \\ & = \log \sum_{\mathbf{z}_D^{\text{new}}} p(\mathcal{D}^{\text{new}}, \mathbf{z}_D^{\text{new}} | \alpha, \eta) = \\ & = \log \sum_{\mathbf{z}_D^{\text{new}}} \int p(\mathcal{D}^{\text{new}}, \mathbf{z}_D^{\text{new}}, \beta_{1:K}, \theta_D^{\text{new}} | \alpha, \eta) d\theta_D^{\text{new}} d\beta_{1:K} = \quad (184) \\ & = \log \sum_{\mathbf{z}_D^{\text{new}}} \int p(\mathcal{D}^{\text{new}} | \mathbf{z}_D^{\text{new}}, \beta_{1:K}) p(\mathbf{z}_D^{\text{new}} | \theta_D^{\text{new}}) \times \\ & \times p(\theta_D^{\text{new}} | \alpha) p(\beta_{1:K} | \eta) d\theta_D^{\text{new}} d\beta_{1:K} \end{aligned}$$

Questa espressione dimostra con chiarezza che la log-verosimiglianza marginale

racchiude l'incertezza sui parametri del modello, poiché essi sono stati integrati dall'espressione complessiva dalla distribuzione di probabilità congiunta dei dati, delle variabili latenti e dei parametri stessi. Naturalmente, come abbiamo già sottolineato nella Sezione (2), la (184) non è trattabile, poiché la somma su tutte le possibili assegnazioni $\mathbf{z}_D^{\text{new}}$ contiene $K^{N_d^{\text{new}}}$ termini per ciascun documento di test. Una soluzione più semplice consiste nel non tener conto dell'incertezza sui parametri, supponendo che essi siano prefissati a dei valori ragionevoli, indichiamoli con $\tilde{\theta}_D$ e $\tilde{\beta}_{1:K}$. Con questa approssimazione, le distribuzioni $p(\theta_D^{\text{new}}|\alpha)$ e $p(\beta_{1:K}|\eta)$ collasano su una massa puntuale, e la (184) si riduce a (Heinrich, 2008; Grün e Hornik, 2011; Papanikolaou et al., 2017):

$$\begin{aligned} & \log p(\mathcal{D}^{\text{new}}|\tilde{\beta}_{1:K}, \tilde{\theta}_D^{\text{new}}) = \\ & = \log \sum_{\mathbf{z}_D^{\text{new}}} \int p(\mathcal{D}^{\text{new}}|\mathbf{z}_D^{\text{new}}, \tilde{\beta}_{1:K}) p(\mathbf{z}_D^{\text{new}}|\tilde{\theta}_D^{\text{new}}). \end{aligned} \quad (185)$$

Possiamo allora ulteriormente elaborare questa espressione, per arrivare ad una espressione di calcolo della log-verosimiglianza, riportata in Heinrich (2008) senza particolari giustificazioni. Osserviamo innanzitutto che:

$$\begin{aligned} & \sum_{\mathbf{z}_D^{\text{new}}} \int p(\mathcal{D}^{\text{new}}|\mathbf{z}_D^{\text{new}}, \tilde{\beta}_{1:K}) p(\mathbf{z}_D^{\text{new}}|\tilde{\theta}_D^{\text{new}}) = \\ & = \sum_{i=1}^K \prod_{d^{\text{new}}=1}^T \prod_{n=1}^{N_d^{\text{new}}} \prod_{j=1}^{|V|} \tilde{\beta}_{ij}^{w_{dn}^{j,\text{new}}} \prod_{s=1}^K \tilde{\theta}_{ds}^{\text{new}, z_{dn}^{s,\text{new}}} = \\ & = \sum_{i=1}^K \prod_{d^{\text{new}}=1}^T \prod_{n=1}^{N_d^{\text{new}}} \prod_{j=1}^{|V|} (\tilde{\beta}_{ij} \tilde{\theta}_{di}^{\text{new}})^{w_{dn}^{j,\text{new}}} = \\ & = \prod_{d^{\text{new}}=1}^T \prod_{n=1}^{N_d^{\text{new}}} \prod_{j=1}^{|V|} \left[\sum_{i=1}^K (\tilde{\beta}_{ij} \tilde{\theta}_{di}^{\text{new}})^{w_{dn}^{j,\text{new}}} \right]. \end{aligned} \quad (186)$$

Poiché l'indicatore w_{dn}^j assume sempre valore 0 oppure 1 su qualsiasi documento (sia esso parte del corpus di training o di quello di test) è ovvio che possiamo scrivere:

$$\begin{aligned}
& \prod_{d^{\text{new}}=1}^T \prod_{n=1}^{N_d^{\text{new}}} \prod_{j=1}^{|V|} \left[\sum_{i=1}^K (\tilde{\beta}_{ij} \tilde{\theta}_{di}^{\text{new}})^{w_{dn}^{j,\text{new}}} \right] = \\
& = \prod_{d^{\text{new}}=1}^T \prod_{n=1}^{N_d^{\text{new}}} \prod_{j=1}^{|V|} \left[\sum_{i=1}^K (\tilde{\beta}_{ij} \tilde{\theta}_{di}^{\text{new}}) \right]^{w_{dn}^{j,\text{new}}}.
\end{aligned} \tag{187}$$

Ma allora, l'espressione approssimata della log-verosimiglianza marginale del corpus di test diventa:

$$\begin{aligned}
& \log p(\mathcal{D}^{\text{new}} | \tilde{\beta}_{1:K}, \tilde{\theta}_D^{\text{new}}) = \\
& = \log \prod_{d^{\text{new}}=1}^T \prod_{n=1}^{N_d^{\text{new}}} \prod_{j=1}^{|V|} \left[\sum_{i=1}^K (\tilde{\beta}_{ij} \tilde{\theta}_{di}^{\text{new}}) \right]^{w_{dn}^{\text{new},j}} = \\
& = \log \prod_{d^{\text{new}}=1}^T \prod_{j=1}^{|V|} \prod_{n=1}^{N_d^{\text{new}}} \left[\sum_{i=1}^K (\tilde{\beta}_{ij} \tilde{\theta}_{di}^{\text{new}}) \right]^{w_{dn}^{\text{new},j}} = \\
& = \log \prod_{d^{\text{new}}=1}^T \prod_{j=1}^{|V|} \left[\sum_{i=1}^K (\tilde{\beta}_{ij} \tilde{\theta}_{di}^{\text{new}}) \right]^{n_{dj}^{\text{new}}} = \\
& = \sum_{d=1}^T \sum_{j=1}^{|V|} n_{dj}^{\text{new}} \log \left[\sum_{i=1}^K (\tilde{\beta}_{ij} \tilde{\theta}_{di}^{\text{new}}) \right].
\end{aligned} \tag{188}$$

Questa espressione può essere utilizzata direttamente nella (182) per approssimare la perplexity: ci resta solo da quali stime possiamo utilizzare concretamente per $\tilde{\beta}_{ij}$ e $\tilde{\theta}_{di}^{\text{new}}$. Per quanto riguarda il parametro globale β_{ij} , l'ipotesi implicita è che i documenti di test provengano dallo stesso processo generatore dei dati di training, e quindi che la distribuzione sui topic (che non dipende dallo specifico documento) sia la stessa tanto sul corpus di training che su quello di test. Poiché l'apprendimento dei topic avviene sul corpus di training, appare ragionevole che il parametro globale β_{ij} sia prefissato al valore stimato durante la fase di apprendimento, ossia $\hat{\beta}_{ij}^G$. In altre parole, è come se i nodi stocastici corrispondenti al parametro matriciale $\beta_{1:K}$ sparissero dalla specificazione del modello quando utilizziamo il corpus di test, poiché questi documenti non vengono utilizzati una seconda volta per ri-apprendere

nuovamente $\beta_{1:K}$.

Per i parametri locali θ_{di}^{new} l'apprendimento può avvenire tramite collapsed Gibbs Sampling: poiché $\beta_{1:K}$ è prefissato al valore stimato sul corpus di training, l'equazione di aggiornamento si semplifica in:

$$p(z_{dn} = i | \mathbf{z}_{-(dn)}, \alpha, \mathcal{D}^{\text{new}}) \propto \hat{\beta}_{ij}^G [C(d, i)_{-(dn)} + \alpha_i]. \quad (189)$$

Una volta ottenute (documento per documento) le stime di ciascun θ_{di}^{new} mediante la (173), possiamo inserirle nella (188) al posto di $\tilde{\theta}_{di}^{\text{new}}$. Si noti che la stima Bayesiana di θ_{di}^{new} può essere basata sulla media di L stime ottenute prendendo L stati della catena di Markov. Infatti, poiché β_{ij} è prefissato, nella fase di test non abbiamo problemi di label switching. Naturalmente, invece di utilizzare un singolo corpus di test, la valutazione della perplexity può essere effettuata all'interno di uno schema di tipo cross-validation.

In modo alternativo ai calcoli approssimati che abbiamo descritto, possiamo cercare di valutare l'integrale (184) mediante quadrature numeriche ovvero algoritmi MCMC specializzati a tal fine. Per questo approccio, sicuramente più corretto da un punto di vista formale ma molto più impegnativo computazionale, rimandiamo a Wallach e al. (2009b) e Buntine (2009).

9. Altri sviluppi più recenti del modello LDA

In questa sezione discuteremo di alcune specializzazioni del modello LDA, proposte successivamente alla pubblicazione del paper originale di Blei et al. (2003). A differenza delle sezioni precedenti l'enfasi non verrà posta sull'analisi dettagliata della struttura dei modelli e dei relativi metodi di stima. Piuttosto, l'approccio sarà quello di una breve rassegna, destinata a puntualizzare quelli che sono stati gli sviluppi più importanti nel campo dei modelli a topic latenti, e a fornire un insieme di puntatori alla letteratura essenziale sull'argomento. Rassegne più estese e comprensive sull'argomento sono Kherwa e Bansal (2018) e Vayansky e Kumar (2020).

Una prima possibile estensione riguarda la natura dell'apprendimento implicato dall'utilizzo del modello LDA standard. Se guardiamo alle distribuzioni di probabilità delle proporzioni dei topic, è evidente che esse rappresentano un potente strumento di riduzione della dimensionalità di ciascun documento in un corpus testuale. Allo stesso modo, le distribuzioni di probabilità sui topic $\beta_{1:K}$ permettono di inferire, senza l'utilizzo di informazione esterna al corpus testuale stesso, la semantica di ciascun documento. Pertanto, il modello LDA si muove all'interno di un'ottica

chiaramente non supervisionata.

In altre situazioni, a ciascun documento può essere associata una variabile di risposta categoriale, oppure una variabile numerica discreta o continua. Ad esempio, tale variabile potrebbe essere rappresentata da una semplice etichetta binaria alla quale è associata una determinata semantica, ovvero dal numero di utenti che hanno trovato un certo documento interessante. In questo caso, l'interesse si sposta in senso supervisionato, e l'obiettivo diventa quello di utilizzare i topic latenti per prevedere la variabile di risposta su nuovi documenti futuri. Questa estensione è esattamente quella che viene tratta all'interno del modello **sLDA (supervised latent Dirichlet allocation)**, così come è stato originariamente proposto in Blei e McAuliffe (2007).

Supponiamo, per semplicità, di considerare il caso in cui a ciascun documento è associata variabile di risposta binaria, ossia $y \in \{0,1\}$. La struttura generativa del modello sLDA è la seguente:

1. generiamo un vettore K -dimensionale:

$$\theta \sim \text{Dirichlet}_K(\alpha).$$

2. per ciascuno degli N termini del documento ($n = 1, 2, \dots, N$) ed indipendentemente l'uno dall'altro:
 - c. generiamo un topic $z_n | \theta \sim \text{Multinomial}_K(\theta)$;
 - d. generiamo un token dalla distribuzione Multinomiale $|V|$ -dimensionale $p(w_n | z_n, \beta)$, cui parametri sono condizionati al topic z_n che è stato generato.
3. generiamo la variabile di risposta y da un modello lineare generalizzato logistico (**GLMs, Generalized Linear Models**; Dobson e Barnett, 2008; Dunn e Smyth, 2018):

$$y | \mathbf{z}, \eta \sim \text{Bernoulli} \left(\frac{\exp(\eta^T \bar{\mathbf{z}})}{1 + \exp(\eta^T \bar{\mathbf{z}})} \right),$$

dove:

$$\bar{\mathbf{z}} = \frac{1}{N} \sum_{n=1}^N z_n.$$

Il significato di $\bar{\mathbf{z}}$ dovrebbe essere chiaro dalla sua stessa definizione: ciascun vettore z_n è un vettore unitario, il cui unico elemento uguale ad 1 indica il topic al quale è stato assegnato il termine in posizione n del documento. Poiché vettori della stessa dimensione si sommano componente per componente, è chiaro che se

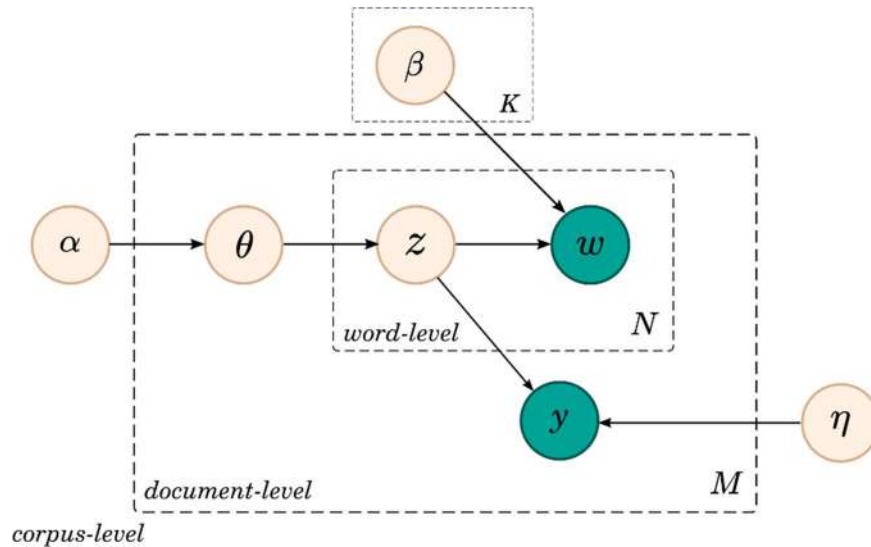
sommiamo z_n rispetto ad n otterremo un vettore K -dimensionale che contiene il numero di termini del documento che sono stati attribuiti a ciascun topic: dividendo per N otterremo la frequenza relativa empirica di occorrenza di ciascun topic nel documento (non direttamente osservabile, poiché le variabili z_n sono latenti). Dovrebbe essere inoltre ben chiaro che \bar{z} è ben distinto da θ , in quanto \bar{z} contiene le frequenze relative effettive di occorrenza di ciascun topic nel documento, mentre θ contiene le frequenze relative attese di occorrenza dei topic, che compaiono nel modello generativo come valore atteso delle componenti marginali della Multinomiale che genera i topic.

Ovviamente, il caso di una risposta binaria non è il solo possibile. In base alla specificazione che abbiamo scritto, la variabile di risposta è un GLM che utilizza come variabili input le frequenze empiriche di occorrenza dei topic (il vettore dei parametri è η). Come è ben noto, la forma distribuzionale della variabile di risposta dei modelli lineari generalizzati è, nel caso più generale possibile, la seguente:

$$p(y|\mathbf{z}, \eta, \delta) = h(y, \delta) \left\{ \frac{(\eta^T \bar{\mathbf{z}})y - A(\eta^T \bar{\mathbf{z}})}{\delta} \right\}, \quad (190)$$

nota come **famiglia a dispersione esponenziale** (Morris, 2006) con parametro di dispersione δ (nel caso della distribuzione Bernoulliana $\delta = 1$). Il lettore familiare con il formalismo dei GLM dovrebbe facilmente riconoscere il significato della (187). Più precisamente: per ogni δ prefissato, la (187) è una **famiglia esponenziale naturale**, per la quale il vettore delle statistiche sufficienti minimali coincide con la funzione identica: si confronti l'espressione precedente con la (41). Inoltre, il valore atteso della variabile di risposta viene espresso come funzione lineare delle variabili di input attraverso l'uso della **funzione di link canonica** (che nel caso della distribuzione Bernoulliana coincide con la funzione logistica).

Figura 10. Rappresentazione grafica del modello sLDA.



L'utilizzo di una famiglia di distribuzioni come la (190) rende possibile considerare in modo unificato tutti i casi particolari di interesse, come ad esempio di quello di una risposta distribuita come una Normale o una Gamma (nel caso continuo), oppure come una Poisson o una Binomiale Negativa (nel caso in cui la risposta sia una variabile di conteggio), ovvero una Binomiale (nel caso in cui la variabile di risposta sia, per ciascun documento, il numero complessivo di successi in un certo numero di ripetizioni di un esperimento Bernoulliano). Dal punto di vista grafico, la rappresentazione del modello sLDA come un grafo orientato è riportata nella Figura 10 alla pagina precedente.

Si noti che la variabile di risposta e il nodo stocastico che descrive i termini del documento hanno un antenato comune, ossia il nodo relativo al topic latente, e quindi non possono essere considerati condizionalmente indipendenti. Dunque, il documento viene generato come un sottoinsieme finito da uno stream infinitamente scambiabile di termini (esattamente come nel modello LDA), e sulla base delle frequenze relative empiriche di occorrenza dei topic viene generata la variabile di risposta. In Blei e McAuliffe (2007) viene fatto notare che una possibile specificazione alternativa potrebbe essere quella di generare le variabili di risposta non sulla base delle frequenze relative empiriche \bar{z} , bensì sulla base delle frequenze relative attese θ .

In questo secondo caso, avremmo l'ovvio vantaggio che tutte le quantità stocastiche coinvolte nel modello risulterebbero infinitamente scambiabili (inclusa la variabile di risposta). Tuttavia, Blei e McAuliffe (2007), sulla base di esperimenti di

simulazione, affermano che la performance previsiva di questo modello supervisionato alternativo risulta essere generalmente inferiore a quello del modello supervisionato ‘standard’ rappresentato nella Figura 10. Ciò in quanto le frequenze relative attese di occorrenza dei topic hanno un potere previsivo, per il singolo documento, generalmente inferiore alle frequenze relative empiriche che sono alla base della generazione di quel particolare documento.

Dal punto di vista dell’inferenza, anche in questo caso abbiamo a disposizione l’approccio EM variazionale (trattando i parametri α , $\beta_{1:K}$ ed η come prefissati ma incogniti). I dettagli sulle relative equazioni di aggiornamento (che ricalcano, sostanzialmente, quelle già viste per il modello LDA), nonché su come utilizzare il modello appreso sul corpus di training per la previsione delle variabili di risposta relative a documenti futuri, sono disponibili con un certo livello di dettaglio in Blei e McAuliffe (2010), oppure in Lakshminarayanan e Raich (2011). Inoltre, vogliamo mettere in evidenza (anche se i dettagli sono fuori dagli obiettivi di questa breve rassegna) che anche il modello sLDA è stato esteso in numerose direzioni, dando origine ad una ricca classe di modelli a topic latenti supervisionati. Si veda, a tal fine, il Capitolo 27 di Murphy (2013), oppure Zhang e Kjellström (2015).

L’altra estensione del modello LDA nei confronti della quale si registra una notevole attività di ricerca è quella nella quale la dimensione temporale viene specificatamente introdotta nella specificazione del modello. Se le date di pubblicazione dei documenti di corpus di documenti si estendessero su lunghi periodi di tempo, il modello LDA classico potrebbe non riflettere con precisione i cambiamenti che si sono verificati nel corso del tempo su come gli argomenti sono stati trattati dagli autori. Ad esempio, un manoscritto riguardante una specifica malattia pubblicato nei primi anni del ‘900 ha sicuramente una struttura e un contenuto molto diverso rispetto ad un documento sulla stessa malattia pubblicato nei primi anni 2000. Al fine di poter contemplare anche tal scenari, è stato introdotto in Blei & Lafferty (2006) il modello **DTM (dynamic topic model)** che permettere di tenere in considerazione l’evoluzione nel tempo del modo in cui un certo argomento viene trattato.

Dal punto di vista della specificazione generativa, il modello DTM mescola la struttura del modello LDA classico con il modello nello spazio degli stati per serie storiche (Koller e Friedman, 2009). Il modello standard presentato nella Sezione 3 viene esteso supponendo che ciascuna riga del parametro β evolva come un modello nello spazio degli stati con rumore Gaussiano comune, ossia:

$$\beta_{t,i}|\beta_{t-1,i} \sim \mathcal{N}(\beta_{t-1,i}, \sigma^2 \mathbb{I}_{|V|}), \quad i = 1, 2, \dots, K \quad (191)$$

dove $\mathbb{I}_{|V|} \in \mathbb{R}^{|V| \times |V|}$ indica la matrice identica.

Naturalmente, questa specificazione Gaussiana sembra incompatibile con il fatto che ciascuna β_i debba essere una distribuzione di probabilità, che compare nella distribuzione Multinomiale che presiede alla generazione dei topic. Per questo motivo, invece che utilizzare direttamente la quantità vettoriale β_i , gli elementi di questo vettore sono trasformati sulla scala dei parametri naturali della distribuzione Multinomiale nel modo seguente:

$$\pi(\beta_{t,i})_j = \frac{\exp(\beta_{t,ij})}{\sum_{j=1}^{|V|} \exp(\beta_{t,ij})}, \quad j = 1, 2, \dots, |V|. \quad (192)$$

Allo stesso modo, l'evoluzione delle proporzioni dei singoli topic è catturata attraverso il seguente semplice modello dinamico:

$$\alpha_t | \alpha_{t-1} \sim \mathcal{N}(\alpha_{t-1}, \delta^2 \mathbb{I}_K). \quad (193)$$

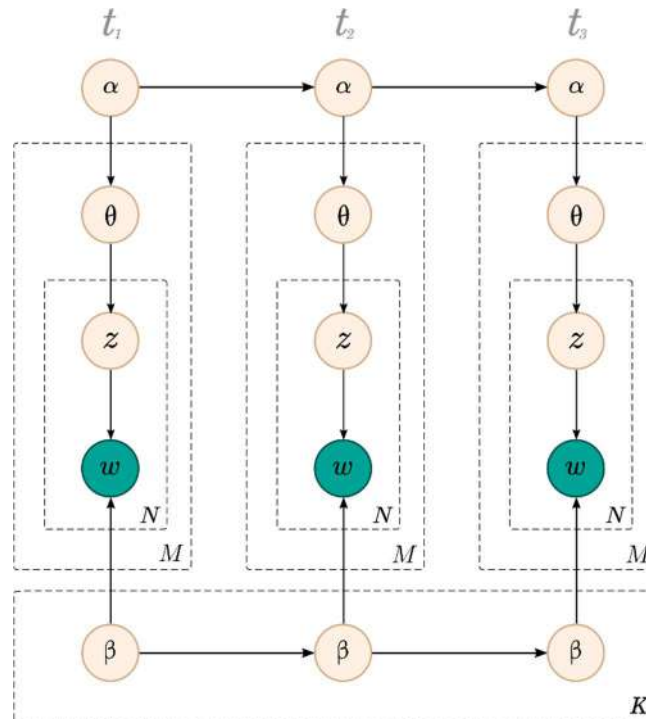
Poiché ciascun vettore α_t può, sulla base della (190), presentare elementi negativi, esso non può essere utilizzato direttamente in una distribuzione di Dirichlet (i cui parametri devono essere tutti positivi). Per questo motivo, generiamo il vettore θ come:

$$\theta \sim \mathcal{N}(\alpha_t, a^2 \mathbb{I}_K). \quad (194)$$

Ma a questo punto, per lo stesso motivo che abbiamo messo in evidenza per $\beta_{t,i}$, anche θ non può essere utilizzato direttamente per le probabilità Multinomiali che sono utilizzate per generare z_n , e pertanto attraverso lo stesso tipo di trasformazione utilizzato nella (189), avremo che (indipendentemente su n):

$$z_n | \theta \sim \text{Multinomial}_K(\pi(\theta)). \quad (195)$$

Figura 11. Rappresentazione grafica del modello DTM.



La rappresentazione grafica del modello DTM è riportata, per tre istanti temporali distinti, nella Figura 11. In essa le frecce orizzontali indicano l'evoluzione dei parametri attraverso successivi istanti temporali, mentre le frecce verticali indicano (come è usuale) la struttura di dipendenza tra i nodi stocastici del modello, per ciascun istante del tempo prefissato. Se rimuoviamo le frecce orizzontali avremo una collezione di modelli LDA standard mutuamente indipendenti: questa specificazione appare alquanto inefficiente, ossia ogni istante del tempo è modellizzato attraverso i suoi propri parametri, che non sono in relazione alcuna con i parametri che descrivono la struttura tematica dei documenti in qualunque altro istante del tempo.

Come tutti gli altri modelli a topic latenti, anche per il modello DTM l'inferenza a posteriori è intrattabile. Blei & Lafferty (2006) introducono un algoritmo variazionale basato sulla teoria del filtro di Kalman. Infine, anche il modello dinamico è stato esteso in varie direzioni: ad esempio, Jähnichen et al. (2018) estendono il modello a tempo continuo introducendo una classe di distribuzione a priori per l'evoluzione dei topic basata sui processi Gaussiani (Rasmussen e Williams, 2005), che può essere

utilizzata non solo in un setting di evoluzione temporale, bensì anche con dati spaziali.

Una ulteriore linea di sviluppo delle ricerche, successive alla pubblicazione del modello LDA originale, ha riguardato lo sviluppo di metodi numerici più efficienti di quelli originariamente proposti. Questa necessità deriva non tanto dalla complessità computazionale degli algoritmi variazionali o di tipo MCMC, bensì piuttosto dall'enorme mole di dati che deve essere spesso trattata quando ci occupiamo di corpus testuali. Con l'algoritmo EM variazionale standard, questo aspetto può dare facilmente luogo a problemi insormontabili di consumo eccessivo di risorse computazionali (sia in termini di memoria che di tempo di occupazione della CPU). Ciò accade in quanto ad ogni step l'algoritmo cicla tra i parametri variazionali locali e globali, e deve pertanto attraversare in ogni caso l'intero corpus ad ogni iterazione, con le conseguenze sopra richiamate quando la dimensione del corpus diventa molto grande.

Per aggirare questo ostacolo, Hoffman et al. (2010) propongono una versione **online** dell'algoritmo standard: riferendoci alla versione presentata per il modello full Bayes della Sezione 7, ad ogni iterazione la versione online ottimizza i parametri variazionali locali relativi a ciascun documento, mentre il parametro variazionale globale η viene aggiornato in modo opportuno utilizzando tutti i documenti che sono stati utilizzati fino a quell'iterazione per ottimizzare i parametri variazionali locali. In questo modo, non dobbiamo attendere che venga attraversato l'intero corpus per poter aggiornare gli iperparametri globali, e l'algoritmo diventa adatto anche a trattare corpus che sono disponibili in streaming (ossia l'informazione disponibile è un flusso continuo distribuito nel tempo).

Anche l'algoritmo collapsed Gibbs Sampling standard è stato oggetto di numerose revisioni e miglioramenti dell'efficienza computazionale complessiva, come ad esempio da parte di Porteous et al. (2008), che propongono uno schema di campionamento MCMC alternativo chiamato **FastLDA**, avente una complessità lineare nel numero di topic K del modello fino ad 8 volte più veloce dell'algoritmo standard (almeno nelle affermazioni degli autori). Altri sviluppi hanno riguardato la parallelizzazione dell'algoritmo su architetture multicore, oppure il calcolo distribuito su cluster di computer collegati in rete; per tutti queste implementazioni particolari rimandiamo alle librerie software, brevemente descritte nella prossima Sezione, nonché alla letteratura specializzata sull'argomento (Newman et al., 2009; Chen et al., 2016).

10. Le principali librerie software disponibili

Abbiamo una vasta disponibilità di software open-source per stimare i parametri del modello LDA e delle sue principali varianti. Qui di seguito ci limitiamo ad elencare le librerie più importanti per le due principali piattaforme di data scienze, ossia Python ed R (R Core Team, 2020).

10.1 Python

La libreria `sklearn` (più estesamente: `Scikit-Learn`) è una delle librerie open-source per l'apprendimento automatico più diffuse su Python (Pedregosa et al., 2011). Essa mette a disposizione il modulo `LatentDirichletAllocation`, che implementa una versione online dell'algoritmo variazionale per l'apprendimento del modello. Tuttavia, è possibile usare anche l'apprendimento batch standard modificando un parametro opportuno. Possiamo scegliere il numero di core da utilizzare, fino ad assorbire tutte le risorse disponibili.

`Tomotopy` (<https://github.com/bab2min/tomotopy>) è a sua volta basata sulla libreria `tomoto` (Topic Modeling Tool) scritta in C++. Nello specifico, è disponibile l'algoritmo collapsed Gibbs Sampling tanto per il modello LDA quanto per le sue principali varianti (come, ad esempio, la versione supervisionata). Come `sklearn`, anche `tomotopy` permette di sfruttare le architetture multicore per velocizzare le iterazioni del campionamento.

`Gensim` (Řehůřek e Sojka, 2010) è una libreria che fornisce funzioni per la stima di modelli a topic latenti, nonché funzioni più specifiche per compiti di ritrovamento automatico dell'informazione, come ad esempio l'associazione di specifici termini di ricerca ad un documento (document indexing), oppure la ricerca di un documento che maggiormente si avvicina al bisogno informativo espresso mediante una query testuale (similarity retrieval). L'implementazione è molto simile a quella di `sklearn` poiché basa sullo stesso script. Anche `gensim` permette di scegliere il numero di core da usare al fine di velocizzare l'apprendimento sfruttando architetture multicore. Una particolarità di `gensim` risiede nel fatto che essa consente di effettuare l'apprendimento distribuito su un cluster di computer.

La libreria `lda` (<https://pypi.org/project/lda/>) apprende i topic latenti usando l'algoritmo collapsed Gibbs sampling, utilizzando una interfaccia simile a quella di `sklearn`. Il codice è scritto sia in C che in Python, con lo svantaggio che non è possibile sfruttare architetture multicore.

Infine, `PyMC3` è una libreria che fornisce uno 'zucchero sintattico' per poter specificare in modo semplice e intuitivo modello Bayesiani gerarchici complessi. Esso

permette di specificare il modello LDA e le sue principali varianti, e dal punto di vista inferenziale implementa una versione approssimata ma particolarmente efficiente dell'algoritmo EM variazionale, nota come ADVI (Automatic Differentiation Variational Inference; Kucukelbir et al., 2017).

10.2 R

La libreria `lda` (<https://github.com/slycoder/R-lda>) implementa l'algoritmo collapsed Gibbs Sampler, ed è scritta interamente in C. Contiene anche funzioni ad hoc per l'analisi della distribuzione a posteriori dei parametri del modello.

`Topicmodels` è una delle librerie più importanti disponibili in questo ambiente (Grün e Hornik, 2011): fornisce gli strumenti base per apprendere topic latenti sia tramite inferenza variazionale che collapsed Gibbs sampling, e si appoggia sulla libreria `tm`, per la gestione delle matrici dati derivanti da corpus testuali (Feinerer et al., 2008). Più specificamente, `topicmodels` fornisce un'interfaccia al codice originale scritto in C da David M. Blei dell'algoritmo EM variazionale sul quale è stata basata la stesura di Blei et al., (2003), nonché al codice dell'algoritmo collapsed Gibbs Sampling scritto in C++ e descritto in Phan et al. (2008).

La libreria `mallet` mette a disposizione una interfaccia per la libreria MALLETT (MACHINE Learning for Language Toolkit) scritta in Java (McCallum, 2002), consentendo di addestrare modelli e caricare i risultati direttamente in R. Mallet fornisce funzioni per l'elaborazione del linguaggio naturale, classificazione supervisionata di documenti testuali, clustering, stima Bayesiana di modelli a topic latenti, e permette anche di gestire altri problemi tipici del trattamento di corpus testuali. L'apprendimento del modello LDA avviene tramite una versione ottimizzata dell'algoritmo collapsed Gibbs Sampling.

Il pacchetto `textmineR` è descritto in Jones (2019), e contiene numerose funzioni per il data mining dei dati testuali e per i modelli a topic latenti, fornendo anche funzioni specifiche per l'analisi esplorativa dei topic appresi, al fine di valutarne la qualità e la coerenza interna. L'apprendimento del modello può avvenire tramite collapsed Gibbs Sampling o Variational EM. Consente l'interoperabilità anche con altre librerie, come ad esempio `mallet`.

10.3 Software di supporto

`R/LDAvis` fornisce una visualizzazione interattiva tramite browser di modelli LDA già addestrati, al fine di aiutare l'utente nell'interpretazione dei topic appresi (Sievert e Shirley, 2014).

Anche `python/tmtoolkit` (<https://pypi.org/project/tmtoolkit/>) fornisce varie funzioni per la valutazione ed esplorazione dei topic appresi. Contiene anche funzioni per interfacciarsi alle librerie `sklearn`, `lda` e `gensim`.

11. Un esempio

Una breve illustrazione delle potenzialità dei metodi esposti è basata sul dataset Coronavirus (COVID-19) Tweets Dataset, che raccoglie a partire dal 20 Marzo 2020 tutti i tweet a livello mondiale che hanno un contenuto collegato alla pandemia di COVID-19 (Lamsal, 2020). Ad oggi (30/10/2020), sono stati raccolti circa 675 milioni di tweet: data l'enorme dimensione del dataset, ci siamo limitati ad utilizzare i tweet ricompresi tra il 20 e il 25 marzo. In questo periodo, la gran parte dei tweet proveniva dall'India, poiché questa Nazione è stata la prima ad implementare il lockdown dopo l'Italia. Inoltre, anche limitandoci a questa finestra, il dataset contiene più di 7 milioni di tweet.

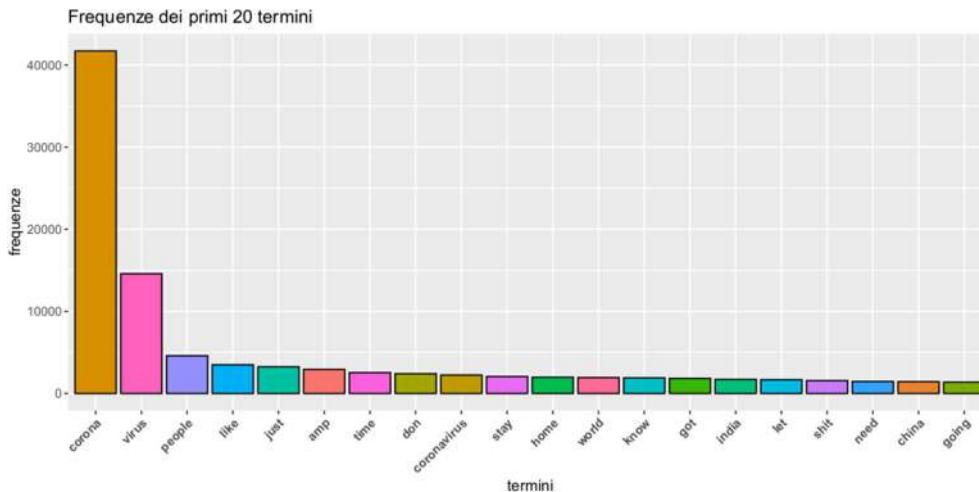
Il dataset contiene, in realtà, solo l'elenco dei tweet ID: questo elenco di ID è stato immesso nel software Hydrator [<https://github.com/docnow/hydrator>], che consente di ottenere la lista di tweet associati (operazione di hydrating). Da questi, sono stati campionati casualmente circa 140K tweet grezzi. Tuttavia, molti tweet entrati in questo set sono risultati essere duplicati (retweet), oppure non in inglese. Dopo le operazioni di pulizia, il dataset finale conteneva circa 55K tweet.

Nella costruzione del modello abbiamo utilizzato dapprima l'algoritmo variational EM così come implementato da `sklearn` in Python. Al fine di effettuare un confronto tra i risultati di un algoritmo variazionale e uno MCMC, i risultati di `Sklearn` sono stati confrontati con quelli basati su collapsed Gibbs Sampling, ottenuti tramite `topicmodels` in R.

La prima fase è quella del pre-processamento dei tweet grezzi: il punto più critico è la tokenizzazione, ossia la suddivisione di ciascun documento nei termini dei quali è composto, e l'eliminazione dei termini irrilevanti per la decodifica semantica dei testi (le cosiddette stopwords, ossia tutti quei token che appaiono con elevata frequenza in tutti i documenti, quali ad esempio le congiunzioni). Al termine di questa fase, i token con le 20 frequenze di occorrenza più elevate sono risultati quelli riportati nella Figura 12. I due token con maggior frequenza di occorrenza sono `corona` e `virus`, come ci si poteva aspettare. Troviamo anche altre parole legate a commenti sull'attualità, visto che il 24 Marzo 2020 in india è stato dichiarato il lockdown. Troviamo anche token relativi alla situazione in altri paesi (`time`, `going`, `world`,

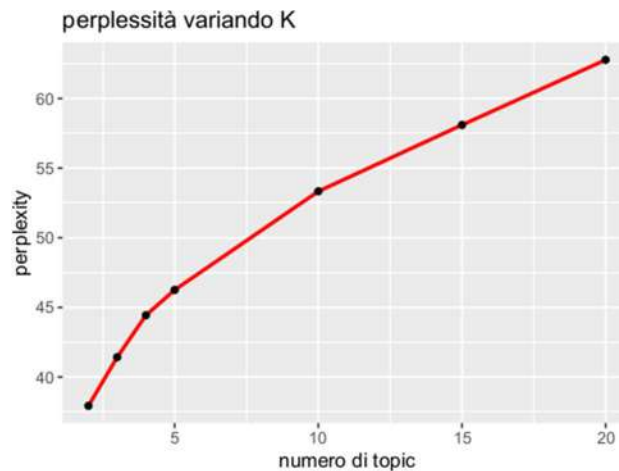
china). Popolari anche i suggerimenti sul restare in casa: stay, home.

Figura 12. I termini più frequenti in un sottoinsieme del dataset Coronavirus (COVID-19) Tweets Dataset, relativo alla finestra compresa tra il 20 e il 25 Marzo 2020.



Per scegliere il numero di topic abbiamo utilizzato una 5-fold cross-validation e la funzione `LatentDirichletAllocation` di `sklearn` (che implementa la versione online dell'algoritmo EM variazionale), con K variabile nell'insieme $\{2, 3, 4, 5, 10, 15, 20\}$. Per ciascun K abbiamo calcolato la perplexity media sui 5 fold. Come sappiamo, il modello ottimale corrisponde a quel valore di K per il quale la perplexità media ha il valore più basso. In questo caso, $K^{\text{opt}} = 2$:

Figura 13. Perplexità del modello LDA, ottenuta mediante 5-fold CV ed utilizzando l’algoritmo variational EM implementato nel modulo LatentDirichletAllocation di sklearn.



Una volta addestrato il modello in corrispondenza di K^{opt} , abbiamo ottenuto i seguenti 30 top-token per ciascuno dei due topic, riportati nella Figura 14. Per quanto riguarda l’interpretazione del primo topic, è evidente che esso riguarda per lo più quei tweet che contengono commenti legati alla politica, all’economia o legati alla religione e spiritualità. Ad esempio, india e narendramodi sono spesso usati dagli utenti che commentano la decisione del governo indiano sul lockdown (Narendra Modi è l’attuale primo ministro indiano), avvenuta il 24 Marzo 2020. Abbiamo anche (Satpal) maharaj, ministro del turismo e della cultura in India. Inoltre, government e country appaiono legati ai commenti sul Governo e sulla gestione dell’emergenza, eventualmente rinforzati anche da pandemic e health. Inoltre, anche god rientra tra i primi 15 token, suggerendo che anche l’argomento religione abbia un peso in questo topic. Tuttavia, spesso può essere usato anche come imprecazione, senza alcun commento religioso. Più in basso troviamo anche economy e financial per i tweet che fanno commenti sull’economia. In questo topic rientrano anche i tweet che riassumono lo stato generale della pandemia (death, news, people) o che invitano a restare in casa: stay e home appaiono tra i primi dieci.

Il secondo topic sembra riguardare la società in generale. Ad esempio, troviamo termini come quarantine, work, usati per discutere del lavoro da remoto. Rientrano in questo topic anche i commenti tweetati dagli utenti per descrivere il modo con cui trascorrono il tempo: ciò è suggerito da time, day, shit e quarantine. Inoltre, trovato spesso il token man, molto utilizzato all’inizio di una risposta. Anche hope e better sono legati a quei commenti di utenti che sperano nel miglioramento della situazione. Stranamente ritroviamo anche trump e chinese, probabilmente usati per fare delle

critiche oppure dei commenti ironici sulla situazione.

Figura 14. I 30 top-token, ossia quelli con la probabilità di occorrenza più elevati, per ciascuno dei due topic, utilizzando il modello addestrato con $K = 2$ e l'algoritmo EM variazionale online implementato in Python/sklearn.

topic_a	topic_b
corona	corona
virus	virus
amp	just
coronavirus	like
people	people
stay	don
home	time
world	got
india	shit
let	think
china	going
fight	know
covid	really
cases	trump
god	quarantine
today	said
19	day
spread	want
help	right
safe	need
pandemic	ve
narendramodi	fuck
covid19	say
country	thing
government	chinese
health	did
news	getting
govt	life
sir	man
ji	make

Figura 15. I 30 top-token, ossia quelli con la probabilità di occorrenza più elevati, per ciascuno dei due topic, utilizzando il modello addestrato con $K = 2$ e il campionamento collapsed Gibbs sampling implementato in topicmodels.

topic1	topic2
corona	corona
virus	virus
like	amp
just	coronavirus
people	stay
time	home
don	world
got	india
shit	let
need	china
going	people
think	fight
day	covid
good	today
really	stop
make	cases
quarantine	god
said	19
want	trump
life	help
right	new
come	spread
say	safe
work	pandemic
ve	know
know	covid19
way	narendramodi
fuck	chinese
thing	country
did	government

Per quanto riguarda il confronto con l’algoritmo collapsed Gibbs sampling, abbiamo utilizzando la libreria topicmodels, effettuando 4000 cicli di burnin e 2000 cicli di campionamento (sempre sul modello ottimale con $K = 2$). L’interpretazione

dei topic produce un risultato simile a quello ottenuto con sklearn, tranne che in questo caso il topic ‘politico’ è il secondo. Notiamo come alcuni termini sono stati correlati al topic più adatto, come trump e cinese, ai quali era stato dato un peso maggiore verso il topic sulla ‘società’ (Figura 15). In genere, l’algoritmo collapsed Gibbs sampling permette di ottenere risultati semanticamente più coerenti, al prezzo di un costo computazionale leggermente più elevato. Ciò implica la necessità di utilizzare l’apprendimento basato sull’approssimazione variazionale quando la dimensione del dataset di training è molto elevata.

12. Discussione e conclusioni

In questo lavoro abbiamo presentato una rassegna ampia dei principali risultati disponibili per il modello LDA, il primo modello a topic latenti presentato in letteratura da Blei et al. (2003). In questo lavoro abbiamo raccolto una collezione di risultati che, almeno finora, erano disponibili in letteratura in modo parziale e frazionato. Per lo stesso motivo abbiamo cercato di raccogliere in un unico contributo tutte le dimostrazioni, al fine di facilitare il compito del ricercatore che volesse approcciarsi al modello LDA, sia per scopi di ricerca pura che applicata.

Il modello LDA standard è sicuramente quello più frequentemente utilizzato per le applicazioni su corpus testuali reali, ma sicuramente non è l’unico disponibile. Anche se è sicuramente adatto a compiti di modellazione generale e di riduzione della dimensionalità, deve essere opportunamente adattato quando l’obiettivo dell’analisi è di natura supervisionata, ovvero la considerazione della dimensione temporale è essenziale per valutare l’evoluzione dei topic del momento.

Tutti questi aspetti sono stati messi brevemente in evidenza nella Sezione 9, ma esistono molte altre limitazioni, che continuano a convogliare interesse di ricerca attorno ai modelli a topic latenti, i quali continueranno verosimilmente ad essere uno degli argomenti trainanti, nei prossimi anni, nel settore del Machine Learning applicato all’elaborazione dei dati testuali.

Senza insistere troppo su problematiche particolari, come ad esempio che il modello LDA standard non è adatto a trattare corpus di documenti testuali troppo brevi, oppure documenti nei quali compaiono più lingue, possiamo senz’altro osservare che anche se versioni molto più espressive e più specializzate del modello di base sono ormai disponibili, la complessità computazionale dei relativi metodi inferenziali è difficilmente sostenibile con grandi volumi di dati (anche con le potenze di calcolo attuali). Per questo motivo, uno degli sviluppi più interessanti degli ultimi anni ha

riguardato la fusione dei modelli a topic latenti con le architetture neurali profonde (deep learning).

Poiché le reti neurali profonde sono approssimatori universali, esse hanno mostrato un elevato potenziale di apprendimento anche per distribuzioni a posteriori complicate altrimenti intrattabili. Nell’inferenza neurale variazionale, la distribuzione a posteriori viene approssimata da una distribuzione variazionale parametrizzata attraverso una rete neurale (Kingma & Welling, 2014; Goodfellow et al., 2016; Miao et al., 2017). In questo modo, tanto il modello generativo a topic latenti, quanto la rete stessa, possono essere addestrati sulla base dell’algoritmo back-propagation (con evidenti vantaggi anche dal punto di vista computazionale).

In conclusione, i modelli a topic latenti costituiscono un importante obiettivo di ricerca, tanto dal punto di vista teorico quanto dal punto di vista applicativo. Sebbene questa introduzione non sia una guida completa, esse può costituire un valido entry-point per chi intenda avvicinarsi a queste tecniche.

Riferimenti bibliografici

- Blei, D. M. (2012). Probabilistic topic models. *Communications of the ACM*, 55(4):77–84. <https://doi.org/10.1145/2133806.2133826>.
- Blei, D. M., Kucukelbir, A., McAuliffe, J. D. (2017). Variational Inference: A Review for Statisticians. *Journal of the American Statistical Association*, 112(518):859–877. <https://doi.org/10.1080/01621459.2017.1285773>.
- Blei, D. M., Lafferty, J. D. (2006). Dynamic topic models. In: *Proceedings of the 23rd international conference on Machine learning - ICML '06*, 113–120. <https://doi.org/10.1145/1143844.1143859>
- Blei, D. M., McAuliffe, J. D. (2007). Supervised Topic Models. In: J.C. Platt, D. Koller, Y. Singer, S.T. Roweis: *Advances in Neural Information Processing Systems 20 (NIPS 2007)*, 121–128.
- Blei, D. M., McAuliffe, J. D. (2010). *Supervised Topic Models*. <http://arxiv.org/abs/1003.0783>
- Blei, D. M., Ng, A. Y., Jordan, M. I. (2003). Latent Dirichlet Allocation. *Journal of Machine Learning Research*, 3:993–1022.
- Buntine, W. (2009). Estimating likelihoods for topic models. In: *Lecture Notes in Computer Science 5828 LNAI*: 51–64. https://doi.org/10.1007/978-3-642-05224-8_6.
- Carlin, B. P., Louis, T. A. (2008). *Bayesian Methods for Data Analysis*. Chapman and Hall/CRC. <https://doi.org/10.1201/b14884>.

- Chen, J., Li, K., Zhu, J., Chen, W. (2016). WarpLDA. In: *Proceedings of the VLDB Endowment*, 9(10):744–755. <https://doi.org/10.14778/2977797.2977801>
- Cover, T. M., Thomas, J. A. (2005). *Elements of Information Theory*. John Wiley & Sons. <https://doi.org/10.1002/047174882X>.
- Dobson, A. J., Barnett, A. (2008). *An Introduction to Generalized Linear Models, Third Edition*. Taylor & Francis.
- Dunn, P. K., Smyth, G. K. (2018). *Generalized Linear Models With Examples in R*. Springer New York. <https://doi.org/10.1007/978-1-4419-0118-7>.
- Feinerer, I., Hornik, K., Meyer, D. (2008). Text Mining Infrastructure in R. *Journal of Statistical Software*, 25(5):1–54. <http://www.jstatsoft.org/v25/i05/>
- Gelman, A., Carlin, J., Stern, H., Dunson, D., Vehtari, A., Rubin, D. (2013). *Bayesian Data Analysis (Third)*. Chapman and Hall/CRC.
- Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press. <https://www.deeplearningbook.org/>.
- Griffiths, T. L., Steyvers, M. (2004). Finding scientific topics. *Proceedings of the National Academy of Sciences*, 101(suppl 1), 5228–5235. <https://doi.org/10.1073/pnas.0307752101>
- Grün, B., & Hornik, K. (2011). topicmodels: An R Package for Fitting Topic Models. *Journal of Statistical Software*, 40:13. <https://doi.org/10.18637/jss.v040.i13>
- Heinrich, G. (2008). *Parameter Estimation for Text Analysis*. <http://www.arbylon.net/publications/text-est.pdf>.
- Hoffman, M. D., Blei, D. M., Bach, F. (2010). Online learning for Latent Dirichlet Allocation. In: J.D. Lafferty, C.K.I. Williams, J. Shawe-Taylor, R.S. Zemel, A. Culotta; *Advances in Neural Information Processing Systems 23, NIPS(2010)*. <https://papers.nips.cc/paper/3902-online-learning-for-latent-dirichlet-allocation>
- Jähnichen, P., Wenzel, F., Kloft, M., & Mandt, S. (2018). Scalable Generalized Dynamic Topic Models. In: *International Conference on Artificial Intelligence and Statistics, AISTATS 2018*, 1427–1435. <http://arxiv.org/abs/1803.07868>.
- Kallenberg, O. (2005). *Probabilistic Symmetries and Invariance Principles*. Springer-Verlag. <https://doi.org/10.1007/0-387-28861-9>.
- Kherwa, P., Bansal, P. (2018). Topic Modeling: A Comprehensive Review. *ICST Transactions on Scalable Information Systems*, 7(24):159623. <https://doi.org/10.4108/eai.13-7-2018.159623>.
- Kingma, D., Welling, M. (2014). Auto-Encoding Variational Bayes. In: *ICLR 2014*. <https://arxiv.org/abs/1312.6114>.
- Koller, D., Friedman, N. (2009). *Probabilistic Graphical Models: Principles and Techniques*. MIT Press.
- Kucukelbir, A., Tran, D., Ranganath, R., Gelman, A., & Blei, D. M. (2017). Automatic Differentiation Variational Inference. *Journal of Machine Learning Research*, 18(1): 430–474.

-
- Lakshminarayanan, B., Raich, R. (2011). Inference in Supervised latent Dirichlet allocation. In: *2011 IEEE International Workshop on Machine Learning for Signal Processing*, 1–6.
- Lamsal, L. (2020). Coronavirus (COVID-19) Tweets Dataset. IEEE Dataport. <http://dx.doi.org/10.21227/781w-ef42>.
- Liu, J. S. (1994). The Collapsed Gibbs Sampler in Bayesian Computations with Applications to a Gene Regulation Problem. *Journal of the American Statistical Association*, 89(427): 958–966.
- Manning, C. D., Raghavan, P., Schütze, H. (2008). *Introduction to Information Retrieval*. Cambridge University Press.
- McCallum, A. K. (2002). *MALLET: A Machine Learning for Language Toolkit*. <http://mallet.cs.umass.edu/>.
- Miao, Y., Grefenstette, E., Blunsom, P. (2017). Discovering Discrete Latent Topics with Neural Variational Inference. <http://arxiv.org/abs/1706.00359>
- Morris, C.N. (2006). Natural Exponential Families. In: *Encyclopedia of Statistical Sciences*. John Wiley & Sons. <https://doi.org/10.1002/0471667196.ess1759.pub2>
- Murphy, K. P. (2013). *Machine Learning: a Probabilistic Perspective*. The MIT Press.
- Newman, D., Asuncion, A., Smyth, P., Welling, M. (2009). Distributed Algorithms for Topic Models. *J of Machine Learning Research.*, 10:1801–1828.
- Jones, T. (2019). *textmineR: Functions for Text Mining and Topic Modeling*. <https://www.rtextminer.com/index.html>.
- Ng, A. Y., Jordan, M. I. (2001). On Discriminative vs. Generative Classifiers: A comparison of logistic regression and naive Bayes. In: T.G. Dietterich, S. Becker, Z. Ghahramani (eds.), *Advances in Neural Information Processing Systems 14 (NIPS 2001)*: 841–848.
- Nigam, K., McCallum, A. K., Thrun, S., Mitchell, T. (2000). Text Classification from Labeled and Unlabeled Documents using EM. *Machine Learning*, 39(2/3):103–134. <https://doi.org/10.1023/A:1007692713085>.
- Papanikolaou, Y., Foulds, J. R., Rubin, T. N., Tsoumakas, G. (2017). Dense Distributions from Sparse Samples: Improved Gibbs Sampling Parameter Estimators for LDA. *Journal of Machine Learning Research*, 18(1):2058–2115.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E. (2011). Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.
- Phan, X. H., Nguyen, L. M., Horiguchi, S. (2008). Learning to classify short and sparse text & web with hidden topics from large-scale data collections. In: *Proceeding of the 17th International Conference on World Wide Web 2008*,

- WWW'08*. <https://doi.org/10.1145/1367497.1367510>.
- Porteous, I., Newman, D., Ihler, A., Asuncion, A., Smyth, P., Welling, M. (2008). Fast collapsed Gibbs sampling for Latent Dirichlet Allocation. In: *Proceeding of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining - KDD 08*, 569–577. <https://doi.org/10.1145/1401890.1401960>.
- R Core Team. (2020). *R: A Language and Environment for Statistical Computing*. <https://www.r-project.org/>.
- Rasmussen, C. E., Williams, C. K. I. (2005). *Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning)*. The MIT Press.
- Řehůřek, R., Sojka, P. (2010). Software Framework for Topic Modelling with Large Corpora. In: *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*. <https://is.muni.cz/publication/884893/en>.
- Robert, C.P., Elvira, V., Tawn, N., Wu, C. (2018). Accelerating MCMC algorithms. *Wiley Interdisciplinary Reviews: Computational Statistics*, 10(5):e1435. <https://doi.org/10.1002/wics.1435>
- Sebastiani, F. (2002). Machine learning in automated text categorization. *ACM Computing Surveys*, 34(1):1–47. <https://doi.org/10.1145/505282.505283>.
- Sievert, C., Shirley, K. E. (2014). LDAvis: A method for visualizing and interpreting topics. *Proceedings of the Workshop on Interactive Language Learning, Visualization, and Interfaces*. <https://github.com/cpsievert/LDAvis>.
- Stephens, M. (2000). Dealing with label switching in mixture models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 62(4): 795–809. <https://doi.org/10.1111/1467-9868.00265>
- Vayansky, I., Kumar, S. A. P. (2020). A review of topic modeling methods. *Information Systems*, 94:101582. <https://doi.org/10.1016/j.is.2020.101582>
- Wallach, H. M., Mimno, D. M., McCallum, A. (2009). Rethinking LDA: Why Priors Matter. In: Y. Bengio, D. Schuurmans, J.D. Lafferty, C.K.I. Williams and A. Culotta (eds.), *Advances in Neural Information Processing Systems 22 (NIPS 2009)*: 1973–1981.
- Wallach, H. M., Murray, I., Salakhutdinov, R., Mimno, D. (2009). Evaluation Methods for Topic Models. In: *Proceedings of the 26th Annual International Conference on Machine Learning*: 1105–1112.
- Zhang, C., Kjellström, H. (2015). How to Supervise Topic Models. In: *ECCV 2014: Computer Vision - ECCV 2014 Workshops*, 500–515. Springer.

Geostatistical analysis of soil reflectance spectra for field-scale digital soil mapping. A case study

Natalia Leone^{1*}, Valeria Ancona¹, Davide Fragnito²,
Domenico Vitale³, Massimo Bilancia⁴

¹Water Research Institute, National Research Council, Viale Francesco de Blasio, 5, Bari,

²Master's Graduate in Statistical, Actuarial and Financial Sciences,

³CMCC Foundation, Euro-Mediterranean Center on Climatic Change, Viterbo,

⁴Ionian Department DJSGE – University of Bari A. Moro, Via Duomo 259, Taranto.

Abstract: Knowledge of field-scale soil variability is essential for sustainable soil management. Traditional techniques, based on soil analysis, are costly and time-consuming. An alternative method would be the use of visible-infrared reflectance spectroscopy coupled with multivariate analysis, specifically principal component analysis (PCA) and geostatistics.

In this study, after brief reviews regarding reflectance spectroscopy, PCA, and geostatistics, we presented a methodological approach for digital soil mapping in a study area of Southern Italy. Reflectance spectra of 240 surface soil samples collected at geo-referenced sites, were decomposed by PCA. The first three components (PC1, PC2, PC3) explained most (98%) of the total variance of the initial data set, therefore, they were considered for the assessment of soil spatial variability by variography and kriging (geostatistics). The resulting PC1, PC2 and PC3 kriging maps were interpreted in the light of the information contents on reflectance spectra and compared with the results of a previous, conventional soil survey. The presented strategy seems to be efficient and reliable for mapping soil spatial variability.

Keywords: Soil reflectance; Principal component analysis (PCA); Geostatistics; Digital Soil Mapping.

* Corresponding author: natalia.leone@ba.irsra.cnr.it

All authors reviewed and revised the manuscript, equally contributing to this work, approved the final version, and agreed to submit the revised manuscript for publication. The authors state that they have no disclosure to declare.

1. Introduction

Soils are rarely homogeneous at different spatial scales (Odlare et al., 2005). Variations in the soil properties, which are particularly evident over a large scale, may occur markedly also within a few hectares of farmland (field scale), due to small changes in topography and thickness of parent material layers or the effects of past human management (Brady and Weil, 2002). Despite this, soils are traditionally treated as homogeneous, with possible adverse effects on crop yield, management costs, and the environment. However, these effects can be contained, if not completely avoided, by adapting soil management to the site's specific conditions, as assessed through the correct knowledge of the within-field variability. This is the purpose of the set of agricultural techniques, better known as "precision agriculture". A way to investigate within-field soil variability could be the production of detailed maps, based on many traditional chemical and physical analyses. However, these analyses are expensive and time-consuming. Therefore, this approach is unsuitable when soils need to be analysed, as in precision farming. Hence, the need to investigate alternative techniques. Recently, particular interest has been shown towards reflectance spectroscopy in the visible and near-infrared visible domain (vis-NIR spectroscopy) (Leone et al., 2012). Vis-NIR reflectance spectroscopy is a rapid, cost-effective, non-invasive, and non-destructive technique that requires only minimal sample preparation and does not require the use of hazardous chemicals (Viscarra Rossel et al., 2006). Vis-NIR reflectance spectroscopy is defined as the ratio between radiation reflected from the surface of a material (the soil, in our case) and the radiation incident on it, at different wavelengths, between 350 and 2500 nm (Drury, 1993).

In the above-mentioned spectral region, each soil constituent has specific absorption properties, due to energy transitions, either electronic (in the visible) or vibrational (in the near-infrared; Leone, 2000a). Therefore, soils with different chemical, physical and mineralogical properties show various spectral features. The latter can be conveniently analysed to acquire either qualitative or quantitative information on these properties (Leone et al., 2012), or to analyse and map the spatial distribution of the soil mantle (digital soil mapping) (Odlare et al., 2005; Viscarra Rossel and Behrens, 2010), in combination with multivariate and geostatistical data analysis. Although promising, the use of vis-NIR spectroscopy, combined with multivariate and geostatistical analysis, has been little used, also due to the lack of knowledge about the basic concepts of vis-NIR spectroscopy, multivariate statistical and geostatistical methods. The present work aims to provide a methodological contribution

to digital soil mapping based on soil spectral reflectance measurements, multivariate statistical methods, and geostatistics.

2. Some basic concepts

2.1 Soil spectral reflectance

Soil is a semi-infinite medium relative to electromagnetic radiation. In other words, electromagnetic radiation incident on the soil is either absorbed within it or is reflected from its surface. The latter can be measured and then related to soil properties (Irons et al., 1989). The fraction of incident flux that is reflected is referred to as spectral reflectance; it is usually measured in the spectral domains of visible (vis, 350-780 nm) and near infrared (NIR, 780-2500 nm). For this, the spectral reflectance in the 350-2500 nm domain is commonly referred to as vis-NIR reflectance. However, sometimes, the spectral domain between 700 and 2500 nm is further divided into near infrared (NIR, 780-1100 nm) and short-wave infrared (1100-2500 nm).

The vis-NIR spectral reflectance of a soil is affected by several chemical and physical properties, referred to as “chromophores”. A given soil sample consists of a variety of chromophores, which vary with environmental conditions. In many cases, the spectral signals related to a given chromophore overlap with those of other chromophores and thereby hinder the assessment of the effect of a given chromophore (Ben-Dor et al., 1999).

Water, organic matter, and minerals are the main chemical chromophores of a soil. Their influence on soil reflectance is related to vibrational motions and electron transitions. The vibrational motions consist of oscillations in the relative positions of bonded atomic cores. The oscillations either stretch molecular bond lengths or bend interbond angles. Energy level transitions involving nuclear vibrations typically result in the absorption or emission of radiation within the infrared portion of the spectrum (Irons et al., 1989). The electronic transitions involve changes in the energy levels of the electrons in soil atoms and molecules. Electronic processes produce absorption bands readily distinguishable from those produced by vibrational processes based on their appearance, and from their general location in the spectrum. These bands occur mostly in the ultraviolet, and extend with diminishing frequency into the visible, but rarely appear in the infrared. The usual limit is an iron band near 1000 nm (Irons et al., 1989; Hunt and Salisbury, 1970). On the other hand, very sharp bands in the near-infrared region are also observed. The frequency of occurrence and intensity of these bands decreases towards the visible range (Hunt and Salisbury,

1970). Below we briefly illustrate the effects of the main chemical chromophores on soil reflectance.

2.1.1 Water

Reflectance spectra of moist soils show prominent absorption bands centred at 1400 and 1900 nm. These bands, along with weaker bands at 970, 1200 nm, and 1777 nm, are attributable to overtones and combinations of fundamental vibrational frequencies of water molecules in the soil. In addition to absorption bands, increasing moisture content generally decreases soil reflectance across the entire reflectance spectrum (Irons et al., 1989).

2.1.2 Organic matter

Organic matter has spectral activity throughout the entire VNIR-SWIR region, especially in the visible region. In general, the spectral reflectance decreases in the entire wavelength range between 400–2500 nm as the organic matter content increases (Hoffer and Johannsen, 1969). Baumgardner et al. (1970, 1985) observed that organic matter plays an important role on soil reflectance when its content exceeds 2% and that the reflectance spectra of soils rich in organic matter often have a concave shape between 500 and 1300 nm, compared with the convex shape of the spectra of soils with low organic matter contents. Due to the strong influence of organic matter in the visible region, a soil becomes darker with increasing organic matter. However, many other soil properties, such as texture, structure, moisture, and mineralogy, can influence this (Hummel et al., 2001), implying that darkness would only be a useful discriminator within a limited geological variation.

Absorption bands by organic in the vis–NIR are often weak and not readily apparent to the naked eye (Stenberg et al., 2010). These bands result from the stretching and bending of NH, CH, and CO groups (Ben-Dor et al., 1999; Bokobza, 1998; Goddu and Delker, 1960).

Bands around 1100, 1600, 1700 to 1800, 2000, and 2200 to 2400 nm have been identified as being particularly important for soil organic carbon (Ben-Dor and Banin, 1995; Dalal and Henry, 1986; Krishnan et al., 1980; Henderson et al., 1992; Morra et al., 1991; Malley et al., 2000; Stenberg et al., 2010). Clark et al. (1990) assigned bands near 2300, 1700, and 1100 nm to combination bands and first and second overtones, respectively, of the C–H stretch fundamentals near 3400 nm.

2.1.3 Minerals

Clay minerals, iron-oxides, and carbonates are the most important minerals affecting

spectral reflectance.

Kaolinite, smectite, and illite are the most abundant clay minerals in soils, particularly those from the Mediterranean region (Torrent, 1995). Kaolinite has characteristic absorption doublets near 2200 and 1400 nm. The absorptions wavelengths near 1400 nm (1393 and 1415 nm) are due to overtones of the O–H stretch vibration near 2778 nm, while those near 2200 nm (2165 and 2207 nm) are attributed to Al–OH bend and O–H stretch combinations.

Smectite has sharp characteristics absorptions bands near 1400, 1900, and 2200 nm. The band near 1400 nm is partly due to the overtone of structural O–H stretching in the octahedral layer of this clay mineral. This band, along with that near 1900 nm, is also attributed to vibrational motions in the water molecules bound in the interlayer lattices, in the form of water adsorbed on particle surfaces and hydrated cations (Bishop et al., 1994). Such water is not present in kaolin. Therefore, the presence of a weak absorption band near 1900 nm may be used as a diagnostic feature for kaolinitic dry soils. Illite also shows absorption bands near 1400, 1900, and 2200, which, however, are much weaker than those of smectite as well as near 2340 and 2445 nm (Post and Noble, 1993). The latter could be used to distinguish between illitic and smectitic soils (Post and Noble, 1993). However, they are weak, and especially the one near 2445 nm may be confused with absorptions due to organic matter. The spectral response of soils is strongly affected by the presence and abundance of iron oxyhydroxides (Leone, 2000a, 2000b, 2011).

Goethite (α -FeOOH) and haematite (α -Fe₂O₃) are, by far, the most common Fe-oxide mineral in soils (Torrent et al., 2007; Zhao et al., 2017). Both these minerals show broad and smooth absorption features in the visible-near infrared region, due to electronic transition. The spectra of goethite exhibit absorption bands in the near-infrared, near 920 nm and four absorption bands in the visible, near 420, 480, 600 nm (Leone, 2011). A band near 1700 nm can also be observed (Zheng et al., 2016). The spectra of haematite are characterised by three main absorption bands near 520, 650 and 880 nm (Viscarra Rossel and Behrens, 2010; Viscarra Rossell et al., 2010; Leone, 2011). The absorptions in the visible region cause the vivid colours of Fe oxides, for example, yellow goethite and red haematite (Stenberg et al., 2010).

Carbonates have several absorption bands in the short-wave infrared, due to overtones and combination bands CO₃ fundamental (Clark et al., 1990). The strongest band occurs near 2335 nm, but some weaker absorptions occur near 2160, 1990 and 1870 nm. In addition to chemical chromophores, as discussed above, the reflectance of light from the soil surface is dependent on several physical chromophores. Among these, soil texture, that is the size distribution of the soil mineral particles, plays a

significant role. As reported in Irons et al. (1989), the reflectance generally increases, and contrasts of absorption features decrease as particle size decreases. This behaviour is characteristic of transparent materials, and most silicate minerals behave transparently in the short-wave region. In contrast, the reflectance of opaque materials decreases as particle size decreases.

In common experience, clayey soils often appear darker than sandy soils even though primary clay particles are much smaller than sandy grains. The difference may be explained in part by the different mineralogies of clay and particles but may also be due to the tendency of clay particles to aggregate.

2.2 *Principal Component Analysis*

For reasons of completeness, we include a brief and modern treatment of principal component analysis (PCA). In what follows the main reference is Hastie et al. (2009), while references more specifically oriented to the type of problem we are dealing with are Odlare et al. (2005) and Viscarra Rossel and Chen (2011).

Suppose that the normalized spectra have been arranged into a rectangular matrix $X \in \mathbb{R}^{N \times p}$, in which N represents the number of sampled spatial locations, while p is the number of wavelengths at which the spectra have been sampled. The data matrix is written in extended form as:

$$X = \begin{pmatrix} x_1^T \\ \vdots \\ x_N^T \end{pmatrix}, \quad (1)$$

which can be interpreted as a collection of N points embedded into a p -dimensional Euclidean space ($p \leq N$). We want to represent these data points through the following reduced rank (rank $q \leq p$) affine model:

$$f(\lambda) = \mu + V_q \lambda, \quad (2)$$

where $\mu \in \mathbb{R}^p$, $V_q \in \mathbb{R}^{p \times q}$ and $\lambda \in \mathbb{R}^q$ is a q -dimensional parameter. In this representation, columns of V_q are assumed to be orthonormal; in other words, they are an

orthonormal basis of $\text{span}(V_q)$. This affine approximation must be optimal with respect to the Euclidean normalization, minimizing the reconstruction error:

$$\min_{\mu, \{\lambda_i\}, V_q} \sum_{i=1}^N \|x_i - \mu - V_q \lambda_i\|^2. \quad (3)$$

It can be easily proved that a partial solution of the above optimization problem is given by the following expression (Hastie et al., 2009):

$$\begin{aligned} \hat{\mu} &= \bar{x}, \\ \hat{\lambda}_i &= V_q^T (x_i - \bar{x}), \end{aligned} \quad (4)$$

where $\bar{x} \in \mathbb{R}^p$ is the p -dimensional vector of column means of the data matrix X . This solution is partial in the sense it depends on the orthogonal matrix V_q ; it can be proved that an explicit expression for V_q can be obtained exploiting the following Singular Value Decomposition (SVD):

$$X = UDV^T. \quad (5)$$

In this decomposition we assume that X has been preliminarily centred with respect to column means, $U \in \mathbb{R}^{N \times p}$ and $V \in \mathbb{R}^{p \times p}$ are matrices with orthonormal columns ($U^T U = I_p = V^T V$), with $I_p \in \mathbb{R}^{p \times p}$ the identity matrix) and $D \in \mathbb{R}^{p \times p}$ is a diagonal matrix of non-negative entries. The columns of V are called the right singular vectors of X and are the eigenvectors of the matrix $X^T X \in \mathbb{R}^{p \times p}$ associated with its non-zero eigenvalues. The columns of U are called the left singular vectors of X and are the eigenvectors of the matrix $XX^T \in \mathbb{R}^{N \times N}$ that correspond to its non-zero eigenvalues. The diagonal elements of matrix D are called the singular values of X and are the non-negative square roots of the (common) non-zero eigenvalues of both matrix $X^T X$ and matrix XX^T . We assume that the diagonal elements of D , are in decreasing order and this uniquely defines the order of the columns of U and A (except for the case of equal singular values). Principal components are the columns of the following matrix:

$$\mathbb{Z} = XV = UD. \quad (6)$$

Columns of the \mathbb{Z} matrix are also known as “scores”, as they represent the coordinates of original spectra re-expressed in a new coordinate system obtained through

a rotation of the Euclidean space. Given the rank $q \leq p$, the solution V_q is given by the first q columns of V , and the columns of XV_q are denoted as PC1, PC2, and so on. Orthogonality of principal components is an immediate consequence of orthogonality of right singular vectors (Golyandina and Korobeynikov, 2014).

Let $S_{\mathbb{Z}}$ denote the sample variance-covariance matrix of \mathbb{Z} . The trace of \mathbb{Z} coincides with the sum v of the variances of the original variables:

$$s = \sum_{i=1}^p s_i = \text{trace}(S_{\mathbb{Z}}), \quad (7)$$

in which s_i denotes the sample variance of the i -th principal components. It is also well known that $s_1 \geq s_2 \geq \dots \geq s_p$; based on this property, if we define the cumulative variance as:

$$s_q = \sum_{i=1}^q s_i, \quad (8)$$

the following percentage measures the quota of variance of original variables explained by the first $q \leq p$ principal components:

$$\frac{s_q}{s} \times 100\%. \quad (9)$$

There are several heuristics to select a proper number q^{opt} of principal components to retain. The most common consists in taking the first q^{opt} right singular vectors to capture at least a fixed quota f of the total variance (e.g.: $f = 0.95$ or $f = 0.98$). Other more formal methods do indeed exist, even though they are not used here; the interested reader is sent back to the literature on the subject (Cadima and Jolliffe, 2001; Jolliffe and Cadima, 2016; Orestes Cerdeira et al., 2020).

A physical interpretation of principal components can be based on right singular vectors (or simply eigenvectors) contained in the matrix V , whose elements are also referred to as ‘loadings’. The i -th principal component is a linear combination of

original variables, with weights taken equal to the loadings of the corresponding eigenvector:

$$Xv_i = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_p) \begin{pmatrix} v_{i1} \\ v_{i2} \\ \vdots \\ v_{ip} \end{pmatrix} = v_{i1}\tilde{x}_1 + v_{i2}\tilde{x}_2 + \dots + v_{ip}\tilde{x}_p, \quad (10)$$

where the data matrix X has been written in terms of its p columns (variables).

By analysing the profile of the loadings associated with each principal component it is often possible to reconstruct their meaning, since loadings represents the weight assumed by a specific wavelength in contributing to the determination of the corresponding principal component. However, sometimes it may not be easy to extract a clear interpretation from loadings, and if additional variables measuring soil composition are available it is preferable to assess the association between the scores of principal components and soil variables. This is the approach followed, for example, by Odlare et al. (2005), but see also Viscarra Rossel and Behrens (2010). Sparse principal components (Zou et al., 2006) is another interesting possibility to improve interpretability, reducing the number of explicitly used variables by artificially setting to zero the loadings having absolute values smaller than a predetermined tolerance.

2.3 Geostatistics

The objective of this section is to introduce a model that describes the spatial variation of each principal component Z_j , $j = 1, 2, \dots, q^{opt}$, that simultaneously allows to predict with minimum mean square prediction error its value at spatial locations that have not sampled. For this purpose, each principal component is considered as a regionalized variable:

$$Z_j \equiv \{Z_j(s_i); i = 1, 2, \dots, N, s \in D \subset \mathbb{R}^2\}, \quad j = 1, 2, \dots, q^{opt}. \quad (11)$$

where D is the study area. In this way, observed principal components are considered as a realization of the random function:

$$\{Z_j(s): s \in D \subset \mathbb{R}^2\}. \quad (12)$$

The model used here to describe the spatial variation is the linear model of

regionalization (Wackernagel, 2013):

$$Z_j(s) = \mu_j(s) + W_j(s) + \varepsilon_j(s), \quad (13)$$

where $\mu_j(\cdot) = E(Z_j(\cdot))$ is a deterministic function of spatial coordinates and represents the large scale expected variation, whereas component $\varepsilon_j(\cdot)$ represents the irreducible error, modelled as a spatial White Noise with null expected value and uncorrelated with $W_j(\cdot)$. The function $W_j(\cdot)$ is an intrinsically stationary random function, whose increments are second-order stationary, for which it is possible to define a semi-variogram:

$$\begin{aligned} \gamma_j(h) &= \frac{1}{2} \text{Var}[W_j(s+h) - W_j(s)] = \\ &= \frac{1}{2} E[(W_j(s+h) - W_j(s))^2]. \end{aligned} \quad (14)$$

The semi-variogram $\gamma_j(h)$ changes as the length and direction of vector h change, but it does not depend on its point of application. If $\gamma_j(h)$ depends on h only through its length $\|h\|$ (Euclidean norm), the random function $W_j(\cdot)$ is said to be isotropic. Properties of theoretical variograms $2\gamma_j(\cdot)$ are well known, particularly those for whom a function of a spatial increment h is a valid theoretical variogram (Gaetan and Guyon, 2010). We will also always assume that the theoretical semi-variogram is continuous at zero in the isotropic case. This assumption is equivalent to assume that the random function $W_j(\cdot)$ is L_2 -continuous (or mean-square continuous). Without insisting on mathematical details, L_2 -continuity allows to treat $W_j(\cdot)$ as a random function modelling the local variation on small scale of soil properties, and guarantees the existence of the spatial correlation function as a convenient tool to characterize the soil microstructure (Cressie, 2015).

The optimal linear predictor $Z(\cdot)$ has known expression, either when $\mu(\cdot) = E(Z(\cdot)) = b_0$ or when it has a non-stationary structure, such as $\mu(\cdot) = b_0 + b_1s_x + b_2s_y$, where $s = (s_x, s_y)$ are spatial coordinates. The empirical counterparts of the optimal linear predictor, under each one of these two scenarios, correspond to ordinary kriging and universal kriging, respectively. The kriging equations for estimating the optimal linear predictor presuppose that the functional form of the semi-variogram is known, except for a finite number of parameters. The theoretical semi-variogram has therefore to be replaced by a consistent estimate, and this fact causes several mathematical difficulties, because we do not have theoretical guarantees that the empirical predictor remains optimal in the sense of mean square prediction error

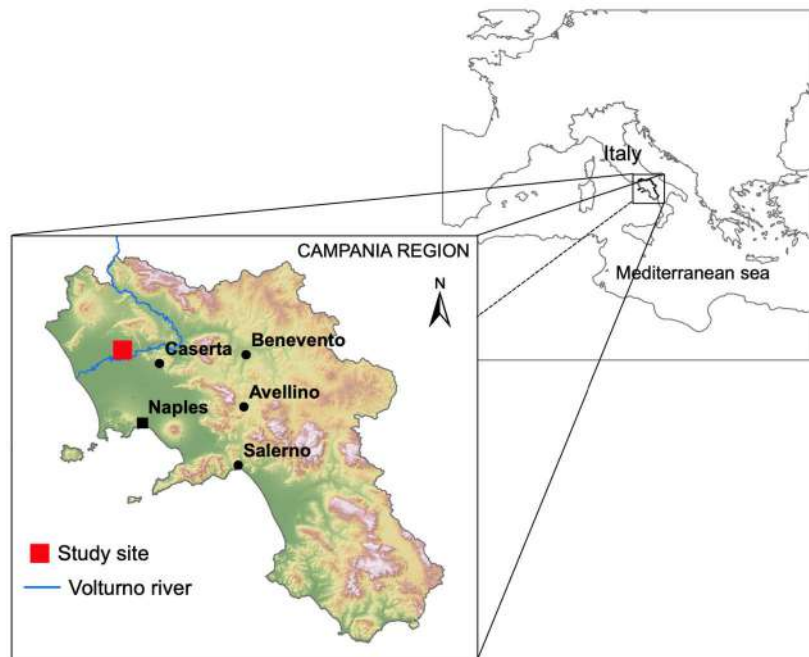
(Bivand et al., 2013). A second difficulty is that kriging has not excellent predictive performances when the likelihood of the data is not Gaussian. In this case it is often convenient to transform the underlying random process, taking for example the logarithm of observed values (as principal components can assume negative values, an offset might eventually be added to ensure that all scores become strictly positive; Varouchakis et al., 2012).

3. Materials and methods

3.1 Site description and sampling

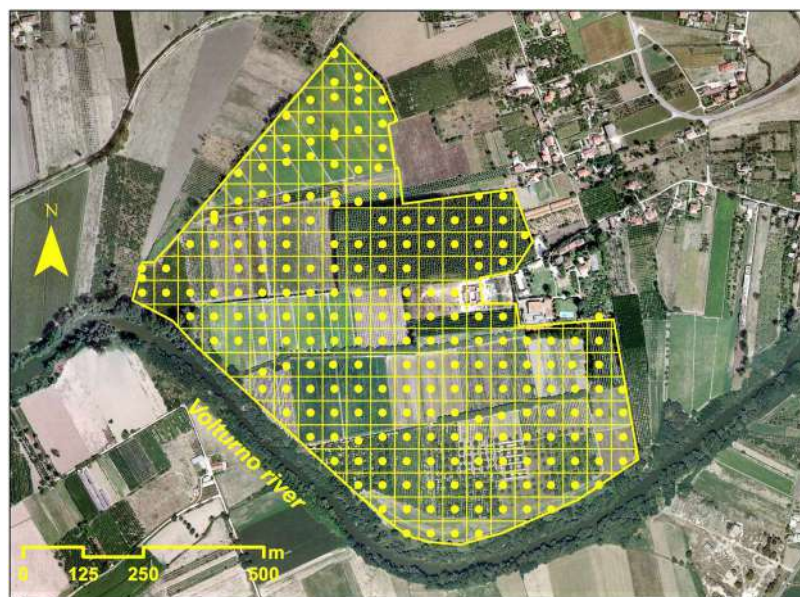
The study covered the entire area (60 ha) of a farm located in the north-western part of the Campania region (Fig. 1), within the municipality of Capua (province of Caserta). This area falls within an abandoned meander of the lower course of the F. Volturno (Aucelli et al., 2014). The production system is mainly oriented towards fruit and cereal growing. The main soil types are Haplic and Fluviic Cambisols and Haplic Luvisols (Grilli et al. 2014; FAO-WRBSR, 2014).

Figure 1. Localisation of the study area (41°06'09" N, 14°11'20" E).



Within the farm, surface soil samples were collected, at a depth of 20-30 cm, at 240 geo-referenced sites (Fig. 2), more or less regularly spaced, falling approximately in the centre of a 50×50 m grid. After collection, the soil samples were transported to the laboratory, air-dried, sieved (2 mm), and finely ground before being subjected to reflectance measurements.

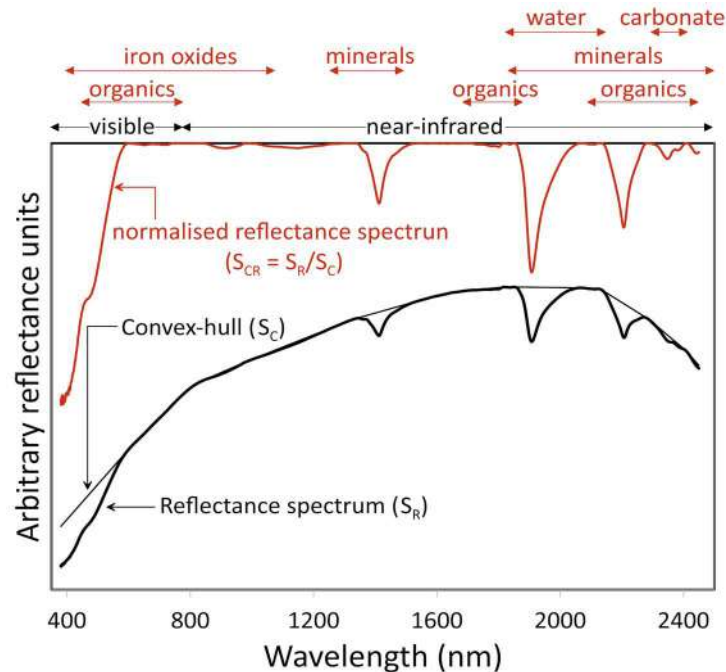
Figure 2. Sampling grid of the investigated field.



3.2 VIS-NIR spectroscopy

The diffuse vis-NIR spectral reflectance was measured in the laboratory, on a residual fraction of soil samples, under controlled light conditions, using the procedure described in Leone et al. (2019). Noisy portions of the measured reflectance spectra, between 350 and 399 nm and between 2451 and 2500 nm, were removed, leaving spectra in the range of 400-2450 nm for the analysis. The resulting reflectance spectra were normalized, using the continuum removal approach (Clark and Roush, 1984), see Fig. 3. To this end, a convex hull was fitted over the original spectral curve, and the absorption spectrum was then calculated considering the ratio between the original reflectance spectrum and the enveloping curve (de Jong, 1992; van der Meer, 1999).

Figure 3. Sample soil vis–NIR spectrum displayed as percent reflectance, convex-hull and continuum removed reflectance. The plot shows regions of the spectrum that hold important information on soil constituents.



3.3 Data pre-processing

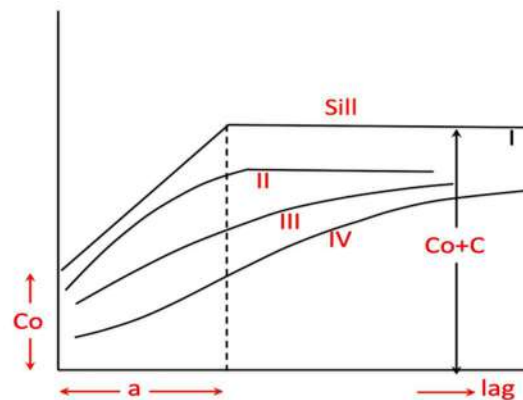
The vis-NIR spectra hold redundant information due to the high degree of correlation between neighbouring wavelengths. For this reason, PCA was performed on the normalized spectra, from which their means (centred data) were subtracted (Viscarra Rossel and Chen, 2011). The initial data were not standardized to the unit of variance since all wavelengths were referred to the same unit, and the differences in their variability were relevant in themselves.

The results of PCA condense the information contained in the spectra. The loadings describe how much each variable contributes to a particular principal component. The PCA reduces the dimensionality of the data, in this case, the reflectance, in a few components, which describe most of the original variance. The first component explains most of the variance, while the subsequent components explain a smaller, progressively decreasing portion.

3.4 Spatial data analysis

The spatial patterns of the first three principal components were analysed using geostatistical analysis. For each principal component, a semi-variogram $\gamma_j(h)$ was estimated which provides a means of quantifying the spatial variation of a variable by measuring the degree of correlation between sampling points separated by a given distance (Webster and Oliver, 2007). The typical parameters on which a theoretical semi-variogram model depends are nugget, range, and sill (Fig. 4).

Figure 4. Examples of semi-variograms: I Linear; II Spherical; III exponential; IV Gaussian. The model parameters nugget (Co), sill ($Co+C$) and range (a) are shown.



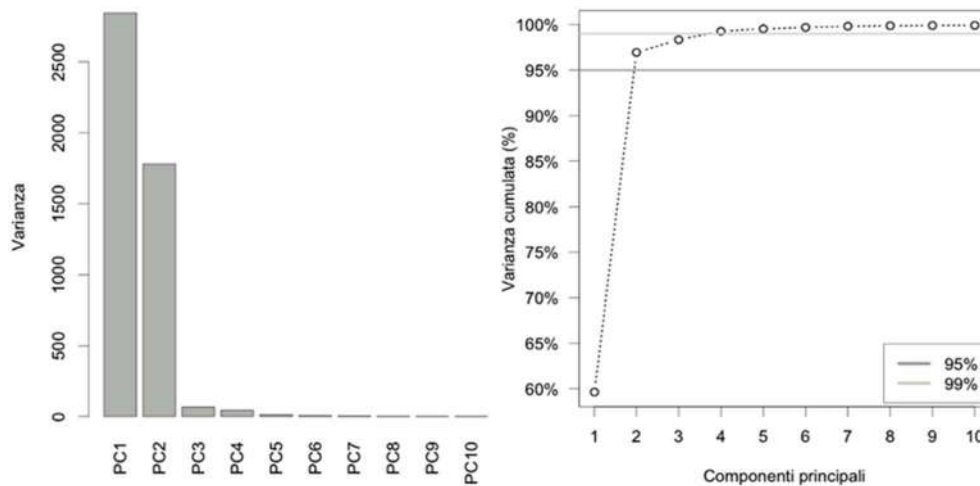
The nugget corresponds to the positive intercept with the y-axis. It is interpretable as the effect of measurement errors, also due to the finite scale at which the phenomenon is observed. Therefore, it increases as the inaccuracy of the measurements increases, i.e., when the sampling interval is too wide. The range is the minimum separation distance at which observations no longer exhibit any spatial correlation. The sill is the point where the theoretical semi-variogram model reaches a limit value (possibly asymptotically) and measures the total variability of the phenomenon. The estimation of the parameters of the theoretical model of semi-variogram is carried out by placing the model in question next to an empirical semi-variogram calculated on data (Cressie, 2015). The estimated theoretical semi-variogram is used to produce a digital map in which the phenomenon under study (in this case, the first three main components analysed separately) is spatially interpolated even outside the sampling sites, through the optimal linear predictor known as kriging.

All the analyses presented next paragraph were carried out using the R 3.6.3 software (R Core Team, 2020).

4. Results and discussions

About 98% of the variance of the original spectral data set is explained by the first three principal components (Fig. 5). In particular, PC1 and PC2 are capable to describe 96.94% of the overall variance, while passing to the first three components together, the overall variance increases to 98.34% (against a modest increase to 99.27% obtained using four components). For this reason, only the scores of the first three principal components were used as input for subsequent analyses.

Figure 5. Variances of the first 10 principal components (a), and percentage of the cumulative variance described by the first 10 principal components. Horizontal bars indicate the number of components which reproduce at least 95% or 99% of the original variability of the normalized vis–NIR spectra (b).



The frequency distribution of the scores of the first three principal components is shown in Fig. 6; from this figure, it is evident that the scores of PC2 present a moderate degree of negative skewness and kurtosis. This reminds us of the probable need for a preliminary logarithmic transformation. It is also evident the presence of some presumable anomalous values (particularly in PC2) located in the left tail of the distribution.

Figure 6. Frequency distribution of scores of the principal components PC1, PC2 and PC3. The empirical coefficient of asymmetry (Skew) and that of kurtosis (Kurt) is also reported above each graph.

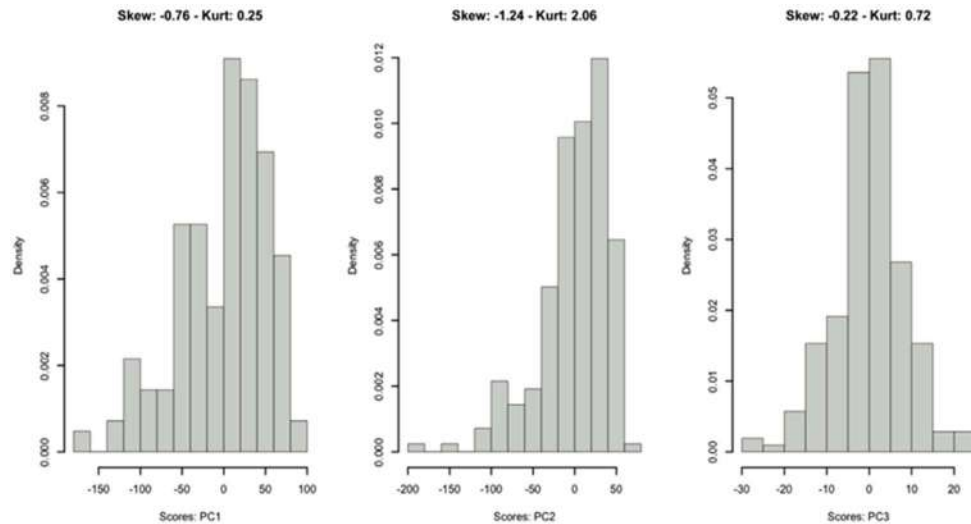
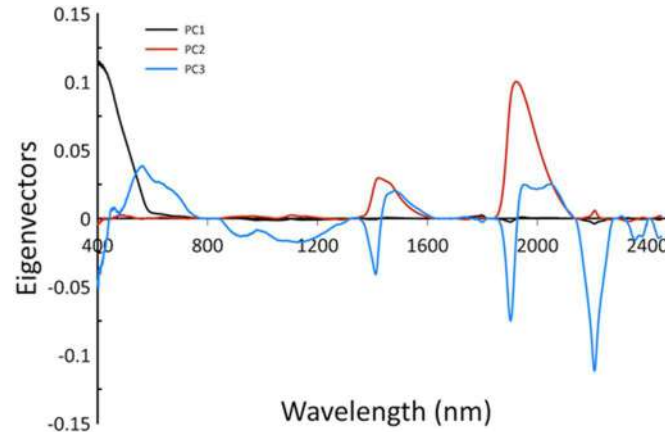


Figure 7 shows the eigenvectors of the first three principal components. The eigenvector of the first principal component showed a steep increase toward the blue and ultraviolet wavelengths, mainly due to a strong iron–oxygen charge transfer band, associated with the presence of iron–oxides, that extend into the ultraviolet (Hunt, 1980). The higher values of loadings in the visible range of the first principal component might be partly due to soil organic carbon (McCauley et al., 1993; Shonk et al., 1991). The eigenvector of the second principal component was dominated by positive loadings near 1400 and 1900 nm, which might be due to 1:1 layer clay minerals (mainly, smectite), specifically to structural O–H stretching mode in its octahedral layer (1400 nm) and combination vibrations of water bound in the interlayer lattices as hydrated cations and water adsorbed on particle surfaces (1400 and 1900 nm) (Bishop et al., 1994; Clark et al., 1990).

Finally, the eigenvector of the third principal component had negative loadings near 1400 and 1900 nm (mainly due to smectite, as previously discussed) and near 2200 nm, due to Al–OH bend in the lattice of 1:2 layer clay minerals (mainly kaolin-ites) (Clark et al., 1990). Table 1 shows the model parameters of the semi-variograms used for the spatialization of the first three principal components.

Figure 7. Eigenvectors (loadings) of the first three principal components (PC₁, PC₂, PC₃).

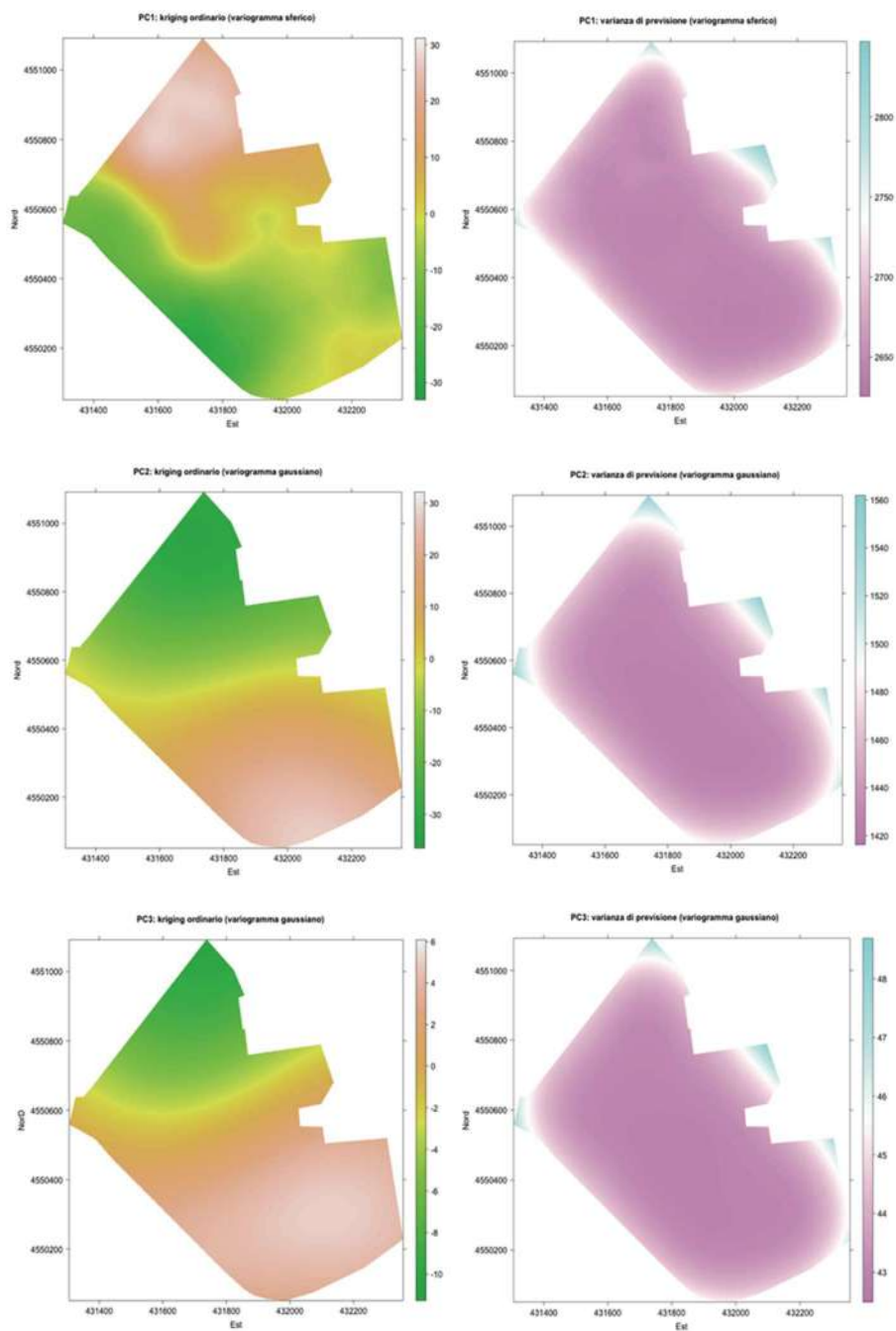
The estimated parameters (using the generalized least squares method; Cressie, 2015) of the theoretical semi-variogram models chosen for the three components were reported in Table 1: the empirical semi-variogram was estimated up to 80% of the maximum theoretically possible distance within the map being studied. The choice of theoretical semi-variogram models, has been largely justified on the basis of an assessment of the goodness of adaptation made ‘by eye’, even though more sophisticated approaches are indeed possible. To map the information content of the main components, we performed the spatial prediction by means of ordinary kriging on a regular discrete grid containing about 2.23×10^6 points.

Table 1. Variogram model parameters of the first three principal components. The Gaussian model reaches its sill asymptotically.

Variable	Model	Nugget	Sill	Range	$C/(C+Co)$
PC1	Spherical	2510	3022	710	0.17
PC2	Gaussian	1400	2390	590	0.41
PC3	Gaussian	42	105	600	0.60

Fig. 8 shows the spectral maps resulting from the spatialization of the scores of the first three principal components. We have reported the optimal linear forecast calculated on the grid (on the left), as well as the corresponding variance of the prediction error (on the right). The quality of the forecasts is somewhat stable across the map, and, as expected, only a modest decline towards the edge of the area under study is highlighted.

Figure 8. On the left: kriging maps of the first three principal components (using variogram models as reported in Table 1). On the right: maps of the relative kriging prediction variance.



The first principal component, as previously discussed, mainly represents the content of iron oxides (Fig. 7). Therefore, we can affirm that the soils of the southernmost area of the investigated company surface (higher PC1 scores), morphologically higher, have higher contents of oxides than iron. This hypothesis is consistent with the geochemical dynamics of soil Fe, strongly influenced by the redox conditions of the medium. In oxidizing conditions, more frequent in the morphologically higher areas of the study area, Fe tends to become insoluble, forming oxyhydroxides. Under reducing conditions, determined by conditions of prolonged water stagnation, more frequent in the morphologically more depressed areas of the study area, the iron compounds dissolve readily, freeing Fe²⁺.

The second principal component represents above all the clay mineral contents of the smectite group. Therefore, the northernmost areas (higher values of the scores) of the study area are likely to be those with the highest clay mineral contents. The third principal component, as already mentioned, is instead inversely related to the clay mineral contents, both of the smectite group and the kaolinite group. This result further confirms the increasing trend of these minerals and, probably, of the finer particle size fractions proceeding from south to north of the company surface. The spectral patterns of the maps relating to the first three principal components (Fig. 8) are consistent with the variability of the soils and their properties, as previously outline in a traditional soil study (Grilli et al., 2014).

5. Conclusions

Vis-NIR spectroscopy, coupled with multivariate statistics and geostatistics, is a useful tool for mapping the spatial variability of soils (digital soil mapping). The reflectance spectra in the Vis-NIR domain contain relevant information on the chemical, physical and mineralogical properties of soils. Multivariate statistical analysis, particularly principal component analysis, is an important tool to condense the highly interrelated reflectance values into a few synthetic variables (principal components), uncorrelated to each other. The geostatistical analysis allows spatializing the principal components and produces spectral maps, which can be interpreted in the light of the known relationships between reflectance and soil properties.

In this study, the spectral maps of the first three principal components have been realized and interpreted. Future studies will be necessary to combine the information contained in the principal component maps, possibly in combinations with other

digital maps (e.g. morphometric maps) using classification and/or data-fusion methods, to produce discretized maps of soil variability, most useful for practical uses.

References

- Aucelli P., Filocamo F., Leone N., Leone A.P. (2014). I paesaggi del Basso Volturno, in Leone, A.P., Buondonno A., Aucelli, P.P.S. (Eds): *Paesaggi e suoli del Basso Volturno per una frutticoltura innovativa*. Iuorio Edizioni, Benevento, pp. 8–40
- Baumgardner, M., Kristof, F.S., Johannsen, C.J., Zachary, A.L. (1970). Effects of organic matter on the multispectral properties of soils. *Proceedings of Indian Academy of Sciences*, 79:413–422
- Baumgardner, M., Silva, L., Biehl, L., Stoner, R. (1985). Reflectance properties of soils. *Advances in Agronomy*, 38:1–44
- Ben-Dor, E., Banin, A. (1995). Near infrared analysis is a rapid method to simultaneously evaluate several soil properties. *Soil Science Society of America Journal*, 59: 364–372
- Ben-Dor, E., Irons, J.R., Epema, G. (1999). Soil reflectance, in Renzc, A. N. (Ed.): *Remote Sensing for the Earth Sciences*. John Wiley & Sons, New York, pp. 111–188
- Bishop, J.L., Pieters, C.M., Edwards, J.O. (1994). Infrared spectroscopic analyses on the nature of water in montmorillonite. *Clays and Clay Minerals*, 42:702–716
- Bivand, R.S., Pebesma, E., Gomez-Rubio, V. (2013). *Applied spatial data analysis with R, 2nd Edition*. Springer New York
- Bokobza, L. (1998). Near infrared spectroscopy. *Journal of Near Infrared Spectroscopy*, 6:3–17
- Brady, N.C., Weil, R.R. (2002) *The Nature and Properties of Soils, 13th Edition*. Pearson Education, Inc., Upper Saddle River, 960
- Cadima, J., Jolliffe, I.T. (2001). Variable selection and the interpretation of principal subspaces. *Journal of Agricultural, Biological, and Environmental Statistics*, 6(1):62–79
- Clark, R.N., Roush, T.L. (1984). Reflectance spectroscopy: quantitative analysis techniques for remote sensing applications. *Journal of Geophysical Research*, 98(B7): 6329–6340
- Clark, R.N., King, T.V.V., Klejwa, M., Swayze, G.A. (1990). High spectral resolution reflectance spectroscopy of minerals. *Journal of Geophysical Research*, 95:12653–12680
- Cressie, N. (2015). *Statistics for Spatial Data*. John Wiley & Sons (Wiley Series in Probability and Statistics)

- Dalal, R.C., Henry, R.J. (1986). Simultaneous determination of moisture, organic carbon and total nitrogen by infrared reflectance spectroscopy. *Soil Science Society of America Journal*, 50:120–123
- de Jong, S.M. (1992). The analysis of spectroscopical data to map soil types and soil crusts of Mediterranean eroded soils. *Soil Technology*, 5(3):199–211
- Drury, S.A. (1993). *Image Interpretation in Geology, 2nd Edition*. Chapman & Hall, London
- FAO-WRBSR (2014). International union of soil science (IUSS) working group world reference base. World reference base for soil resources. In: *International Soil Classification System for Naming Soils and Creating Legends for Soil Maps*. FAO, Rome, Italy, pp. 192
- Gaetan, C., Guyon, X. (2010). *Spatial Statistics and Modeling*, Springer-Verlag New York
- Goddu, R.F., Delker, D.A. (1960). Spectra-structure correlations for the near-infrared region. *Analytical Chemistry*, 32:140–141
- Golyandina, N., Korobeynikov, A. (2014). Basic Singular Spectrum Analysis and forecasting with R. *Computational Statistics & Data Analysis*, 71:934–954
- Grilli, E., Leone, A.P., Buondonno, A. (2014). I suoli dell'Azienda GiòSole, in Leone, A.P., Buondonno, A. and Aucelli, P.P.C. (Eds.): *Paesaggi e Suoli del Basso Volturno per una Frutticoltura Innovativa*. Grafiche Iuorio, Benevento.
- Hastie, T., Tibshirani, R., Friedman, J. (2009). *The Elements of Statistical Learning, 2nd Edition*. Springer New York
- Henderson, T.L., Baumgardner, M.F., Franzmeier, D.P., Scott, D.E., Coster, D.C. (1992). High dimensional reflectance analysis of soil organic matter, *Soil Science Society of America Journal*, 56:865–872
- Hoffer, R.M., Johannsen, C.J. (1969). Ecological potentials in spectral signature analysis, in Johnson, P.L. (Ed.): *Remote Sensing in Ecology*. University of Georgia Press, Athens, pp. 1–29
- Hummel, J.W., Sudduth, K.A., Hollinger, S.E. (2001) Soil moisture and organic matter prediction of surface and subsurface soils using an NIR soil sensor. *Computers and Electronics in Agriculture*, 32:149–165
- Hunt, G.R., Salisbury, J. (1970). Visible and near-infrared spectra of minerals and rocks. Silicate minerals. *Modern Geology*, pp. 283–300
- Hunt, G.R., (1980) - Electromagnetic radiation: the communications link in remote sensing, in Siegal, B.S. and Gillespie, A.R. (Eds.): *Remote Sensing in Geology*. John Wiley & Sons, New York, pp. 5–45
- Irons, J., Weismiller, R., Petersen, G. (1989). Soil reflectance, in Asrar, G. (Ed.): *Theory and Application of optical remote sensing*. Wiley, New York, pp. 66–106

- Jolliffe, I.T., Cadima, J. (2016). Principal component analysis: a review and recent developments. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*. 374(265):2015.0202
- Krishnan, P., Alexander, D.J., Butler, B., Hummel, J.W. (1980). Reflectance technique for predicting soil organic matter, *Soil Science Society of American Journal*. 44:1282-1285
- Leone, A.P. (2000a). Spettrometria e valutazione della riflettanza spettrale dei suoli nel dominio ottico 400 – 2500 nm. *Italian Society of Remote Sensing*, 19:1-26
- Leone, A.P. (2000b). 'Bi-directional reflectance spectroscopy of Fe-oxides minerals in Mediterranean Terra Rossa soils: a methodological approach'. *Agricoltura Mediterranea*. 130:144–154
- Leone, N. (2011) *Uso della spettrometria VNIR e della diffrattometria Rx per la caratterizzazione delle forme pedologiche del ferro: analisi comparativa su modelli di sintesi*, unpublished Bachelor Degree thesis, University of Sannio, Benevento, Italy.
- Leone, A.P., Viscarra-Rossel, A.R., Amenta, P., Buondonno, A. (2012). Prediction of Soil Properties with PLSR and vis-NIR Spectroscopy: Application to Mediterranean Soils from Southern Italy. *Current Analytical Chemistry*. 8(2):283–299
- Leone, A.P., Leone, G., Leone, N., Galeone, C., Grilli, E., Orefice, N., Ancona, V. (2019). Capability of Diffuse Reflectance Spectroscopy to Predict Soil Water Retention and Related Soil Properties in an Irrigated Lowland District of Southern Italy. *Water*, 11(8):1712
- Malley, D. F., Martin, P. D., McClintock, L. M., Yesmin L., Eilers, R. G., Haluschak P. (2000). Feasibility of analysing archived Canadian prairie agricultural soils by near infrared reflectance spectroscopy, in Davies, A.M.C. and Giangiacomo, R. (Eds.): *Near Infrared Spectroscopy: Proceedings of the 9th International Conference*. NIR Publications, Chichester, UK, pp. 579–585
- McCauley, J.D., Engel, B.A., Scudder, C.E., Morgan, M.T., Elliot, P.W. (1993). Assessing the spatial variability of organic matter. *St. Joseph: American Society of Agricultural Engineers*, ASAE Paper No. 93–1555
- Morra, M.J., Hall, M.H., Freeborn, L.L., (1991). Carbon and nitrogen analysis of soil fractions using near-infrared reflectance spectroscopy. *Soil Science Society of American Journal*, 55:288–291
- Odlare, M., Svensson, K. Pell, M. (2005). Near Infrared Reflectance Spectroscopy for Assessment of Spatial Soil Variation in an Agricultural Field. *Geoderma*, 126(3-4):193–202
- Orestes Cerdeira, J., Duarte Silva, P., Cadima, J., Minhoto, M. (2020). subselect: Selecting Variable Subsets. Retrieved from <https://cran.r-project.org/package=subselect>
- Post, J.L., Noble, P.N. (1993). The near-infrared combination band frequencies of dioctahedral smectites, micas, and illites. *Clays and Clay Minerals*, 41: 639–644

- R Core Team (2020). *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria. <https://www.R-project.org/>
- Shonk, G.A., Gaultney, L.D., Schulze, D.G., Van Scoyoc, G.E. (1991). Spectroscopic sensing of soil organic matter content. *Transactions of the American Society of Agricultural Engineers*, 34:1978–1984
- Stenberg, B., Viscarra Rossel, R.A., Mouazen, A.M., Wetterlind, J. (2010). Visible and near infrared spectroscopy in soil science. *Advances in Agronomy*, 107:163–215
- Torrent, J. (1995) *Genesis and properties of the soil of the Mediterranean regions*–University of Naples Federico II, Department of Agro-Chemical Sciences. Arti Grafiche Licenziato, Naples, p. 111
- Torrent J., Liu Q., Bloemendal, J., Barron V. (2007). Magnetic enhancement and iron oxides in the upper luochuan Loess-paleosol sequence, Chinese Loess plateau. *Soil Science Society of America Journal*, 71(5):1570–1578
- van der Meer, F. (1999). Can we map swelling clays with remote sensing? *International Journal of Applied Earth Observation and Geoinformation*. 1(1):27–35
- Varouchakis, E. A., Hristopoulos, D.T., Karatzas, G.P. (2012). Improving kriging of groundwater level data using nonlinear normalizing transformations – a field application. *Hydrological Sciences Journal*, 57(7):1404–1419
- Viscarra Rossel, R.A., Walvoort, D.J.J., McBratney, A. B., Janik, L. J., Skjemstad, J.O. (2006). Visible, near infrared, mid infrared or combined diffuse reflectance spectroscopy for simultaneous assessment of various soil properties. *Geoderma*, 131(1–2): 59–75.
- Viscarra Rossel, R.A., Behrens, T. (2010). Using data mining to model and interpret soil diffuse reflectance spectra. *Geoderma*, 158(1):46–54
- Viscarra Rossel, R.A., Rizzo, R., Dematte, J.A.M. and Behrens, T. (2010). Spatial modelling of a soil fertility index using vis–NIR spectra and terrain attributes. *Soil Science Society of America Journal*, 74:293–1300
- Viscarra Rossel, R. A., Chen, C. (2011). Digitally Mapping the Information Content of Visible–near Infrared Spectra of Surficial Australian Soils. *Remote Sensing of Environment*, 115(6):1443–1455
- Wackernagel, H. (2013). *Multivariate Geostatistics: An Introduction with Applications, 3rd edition*. Springer, Berlin Heidelberg
- Webster, R. and Oliver, M.A. (2007). *Geostatistics for Environmental Scientists, 2nd Edition*. John Wiley & Sons, New York
- Zhao, L., Hong, H., Fang Q., Yin, K., Wang, C., Li, Z., Torrent, J., Cheng, F., Algeo, T.J. (2017). Monsoonal climate evolution in southern China since 1.2 Ma: new constraints from Fe-oxide records in red earth sediments from the Shengli section, Chengdu Basin. *Palaeogeography, Palaeoclimatology, Palaeoecology*, 473:1–15

- Zheng, G., Jiao, C., Zhou, S., Shang, G. (2016). Analysis of soil chronosequence studies using reflectance spectroscopy. *International Journal of Remote Sensing*, 37(8): 1888–1901
- Zou, H., Hastie, T., Tibshirani, R. (2006). Sparse Principal Component Analysis. *Journal of Computational and Graphical Statistics*, 15(3):265–286

Modelli generativi per la sentiment analysis

Claudia Marin*, Fabio Manca.

Dipartimento di Scienze della Formazione, Psicologia, Comunicazione

Riassunto: Il seguente contributo ha lo scopo di fornire un'introduzione e una disamina sistematica dei modelli generativi della sentiment analysis, che rappresentano un'importante famiglia di metodi, per la maggior parte non-supervisionati, che possono essere applicati a tutti i dati testuali grazie alla loro generalità e robustezza. Sono particolarmente utili nell'inferenza per variabili latenti su opinioni e possono eseguire in modo molto efficace analisi congiunte di dati testuali e valutazioni numeriche associate. Consentono, inoltre, la costruzione del lessico dei sentimenti e di lessici specifici per argomenti, l'individuazione di un modello che prenda in considerazione variabili come il tempo, la posizione e le fonti e l'analisi delle preferenze latenti degli opinionisti. Scoprire modelli di opinioni latenti da grandi quantità di dati rende i modelli generativi strumenti essenziali per costruire sistemi intelligenti per la comprensione delle opinioni e per la ricerca nelle scienze sociali computazionali.

Keywords: Sentiment analysis; modelli generativi, big data; social network.

1. Introduzione

I Big Data sono molto voluminosi, hanno una velocità di crescita estremamente rapida, hanno una natura varia, con forme strutturate, semi-strutturate, quasi-strutturate e non strutturate, disponibili in una varietà di formati come documenti, video, audio, immagini, e-mail e "social networking feeds". A causa di questa varietà di fonti, di questa velocità e della veridicità, la qualità non è sempre affidabile e differisce enormemente.

* Autore corrispondente: claudia.marin@uniba.it

Il lavoro qui descritto è frutto di un progetto comune, ma Claudia Marin ha provveduto alla redazione dei paragrafi 3, 4 e 6, mentre Fabio Manca ha provveduto alla redazione dei paragrafi 1, 2 e 5.

Fino a 20 anni fa le organizzazioni non avevano la possibilità di memorizzare in formato elettronico grandi quantità di dati linguistici, come discorsi, lettere, pagine web, articoli scientifici e altri tipi di informazioni verbali. Oggi invece le organizzazioni sono piene di dati non strutturati, utilissimi per migliorare la qualità dei modelli predittivi disponibili.

2. Il text mining

Il text mining è la procedura per ottenere informazioni, estrapolate dal linguaggio naturale, che vengono sintetizzate in un formato strutturato che successivamente viene utilizzato per scopi analitici. E' ormai uno strumento utilizzato in svariati campi disciplinari, quali, ad esempio, la politica per analizzare e interpretare le reazioni del pubblico agli annunci politici, il marketing, i cui reparti sono interessati a valutare la clientela tramite le opinioni sui servizi o le reazioni al lancio di nuovi prodotti, i servizi di sicurezza e le forze dell'ordine per monitorare le liste di sospettati mediante email, messaggi di testo, telefonate e siti web visionati.

Quando si tratta di analisi predittiva, gli stessi strumenti di text mining utilizzati per identificare punti di vista o sentimenti espressi nei documenti possono essere utilizzati anche per estrarre i dati predittori rilevanti e rappresentarli in un formato ben strutturato. Le procedure che portano alla creazione di modelli predittivi possono quindi utilizzare i dati strutturati che sono stati estratti. Generalmente le variabili predittive possono essere o *within document predictors*, atte ad indicare se il testo esprime sentimenti positivi o negativi sull'argomento in discussione, oppure *across document predictors*, che rilevano la somiglianza con altri documenti e sono molto utili, ad esempio, nell'identificare le e-mail fraudolente.

Usare indicatori che rappresentano la frequenza con cui le singole parole o frasi compaiono all'interno di ogni documento è un buon modo per creare variabili predittive, ma per prima cosa è fondamentale procedere alla pulizia dei dati, che è anche definita come "data normalization", ovvero normalizzazione dei dati. Essa implica la rimozione delle stop-words, che sono parole che, data la loro frequenza, sono generalmente ritenute poco significative, lo stemming, che è il processo di riduzione della forma flessa di una parola alla sua radice, detta tema e la standardizzazione, cioè uniformare le parole che hanno lo stesso significato.

Con la normalizzazione viene ridotto il numero delle parole, senza perdere l'essenza della forma. Dopo questa operazione possono seguire il tagging grammaticale, la lemmatizzazione e la categorizzazione semantica (Bolasco *et al.*, 2004).

3. Sentiment analysis

Uno sviluppo del text mining è dato dalla Sentiment Analysis, che è lo studio computazionale di pensieri, opinioni, sentimenti, valutazioni, atteggiamenti, impressioni, stati d'animo ed emozioni delle persone. Approvazione, ammirazione, odio, pietà, disprezzo sono, ad esempio, sentimenti che andrebbero identificati e utilizzati per generare variabili predittive. È una delle aree di ricerca più attiva nell'elaborazione del linguaggio naturale e le sue applicazioni si sono diffuse anche nelle scienze economico-gestionali e sociali.

Esistono numerosi approcci per eseguire l'analisi del sentiment, ma fondamentalmente è possibile distinguerli in due famiglie: l'approccio basato sul lessico (lexicon-based approach) e l'approccio basato sull'apprendimento automatico (machine learning approach).

L'approccio lexicon-based rientra tra le tecniche di classificazione non supervisionate e consiste nel classificare le parole valutando l'orientamento semantico di frasi e documenti a partire dalla polarità delle parole. All'interno del lessico a ciascuna parola viene assegnata una valutazione positiva o negativa. La valutazione globale del sentimento, che si ottiene unendo le parole che hanno la stessa accezione, viene quindi utilizzata come variabile preliminare all'interno del processo di analisi predittiva. Questo approccio, a sua volta, prevede due metodologie: l'approccio basato sui dizionari (dictionary-based approach) e l'approccio basato su corpora (corpus-based approach).

Gli approcci basati sull'apprendimento automatico si affidano ad algoritmi di intelligenza artificiale (Russell e Norving, 2005) per risolvere i problemi di sentiment analysis. Anche questi approcci si dividono in due categorie: con apprendimento supervisionato (supervised learning) e con apprendimento non supervisionato (unsupervised learning). Nel primo caso il sistema acquisisce conoscenze ed esperienze di classificazione grazie ad un training set di testi già etichettati e classificati, mentre nel secondo caso l'acquisizione di conoscenze ed esperienze per la classificazione deriva dall'estrazione di caratteristiche comuni da un training set di testi non etichettati e classificati.

3.1 *Modelli generativi e modelli discriminativi*

Tra gli approcci basati sull'apprendimento automatico è utile distinguere i modelli generativi dai modelli discriminativi (Bishop, 2006). I modelli generativi sono modelli probabilistici di variabili osservate e non osservate che tentano di analizzare la probabilità congiunta di tutti i dati rilevanti con parametri che possano essere

interpretati come strutture o proprietà latenti nei dati. Adattando un tale modello ai dati osservati, si può ottenere una stima di questi parametri e analizzare le strutture latenti ricavate.

Al contrario, i modelli discriminativi, come il support vector machine (SVM), modellano direttamente i confini decisionali, fornendo un modello solo per le variabili oggetto condizionate alle variabili osservate.

Uno dei maggiori vantaggi dei modelli generativi rispetto ai modelli discriminativi è la capacità di esprimere relazioni complesse tra le variabili osservate e quelle oggetto, anche quando tali relazioni non sono direttamente osservabili. Questa proprietà è di particolare importanza, quando si formalizza la sottile dipendenza tra il sentiment e il contenuto del documento di testo per modellare in maniera più accurata le opinioni.

Nell'ultimo decennio Wang *et al.* (2010, 2011), Moghaddam and Ester (2011) e McAuley and Leskovec (2013) hanno ottenuto promettenti progressi nell'esplorazione dei modelli generativi per la sentiment analysis, in quanto hanno sviluppato modelli che, utilizzando sia i dati testuali che le valutazioni numeriche del sentiment, consentissero un'analisi più approfondita delle opinioni per ottenere non solo sentimenti nascosti, ma anche pesi latenti relativi agli argomenti secondari.

4. Modelli linguistici per il testo

Il modello generativo più semplice per la formazione dei dati di testo è quello di linguaggio N-gram, introdotto per la prima volta nel riconoscimento vocale per distinguere tra parole e frasi che suonano simili (Kats, 1987; Rabiner and Juang, 1993) e successivamente introdotto al recupero delle informazioni per la corrispondenza tra query di parole chiave e documenti di testo (Pont and Croft, 1998; Hiemstra and Kraaij, 1998; Zhai and Lafferty, 2001).

Un modello di linguaggio statistico specifica una distribuzione di probabilità sulle sequenze di parole. Ad esempio, con un modello linguistico stimato su una raccolta di articoli di ricerca scientifica si possono fare affermazioni statistiche su quale sequenza di testo è più probabile che sia generata da un esperto informatico.

Formalmente un modello di linguaggio $P(\varphi_1, \varphi_2, \dots, \varphi_n)$ specifica la probabilità congiunta di una sequenza di parole $\varphi_1, \varphi_2, \dots, \varphi_n$. Usando la regola della catena diventa

$$\begin{aligned} P(\varphi_1, \varphi_2, \dots, \varphi_n) &= P(\varphi_1)P(\varphi_2|\varphi_1)P(\varphi_3|\varphi_1\varphi_2) \dots P(\varphi_n|\varphi_1\varphi_2 \dots \varphi_{n-1}) = \\ &= \prod_{k=1}^n P(\varphi_k|\varphi_1, \dots, \varphi_{k-1}) \end{aligned} \quad (1)$$

dove $P(\varphi_k|\varphi_1, \dots, \varphi_{k-1})$ è una distribuzione multinomiale sulle parole nel vocabolario, data la sequenza di parole $\varphi_1, \dots, \varphi_{k-1}$.

La regola della catena mostra il collegamento tra il calcolo della probabilità congiunta di una sequenza di parole e il calcolo della probabilità condizionata di una parola, date tutte le parole precedenti. Intuitivamente l'equazione (1) definisce il processo di generazione di una sequenza di parole, selezionando ripetutamente la parola successiva rispetto a tutte le parole che vengono prima, fino a raggiungere la lunghezza della sequenza predefinita. Per questo motivo tale modello è definito spesso modello generativo. Purtroppo però la complessità computazionale non risulta ridotta e non vi è alcun modo efficiente per calcolare la probabilità esatta di una parola, data una lunga sequenza di parole precedenti.

I modelli N-gram forniscono una soluzione pratica a questa complessa sfida computazionale, approssimando la sequenza precedente ad un numero finito di parole precedenti, cioè $P(\varphi_k|\varphi_1, \dots, \varphi_{k-1}) = P(\varphi_k)$; di conseguenza,

$$P(\varphi_1, \varphi_2, \dots, \varphi_n) = \prod_{k=1}^n P(\varphi_k). \quad (2)$$

In letteratura, il modello unigramma è anche indicato come modello della borsa di parole (bags-of-words model) nel quale l'ordine delle parole è totalmente ignorato.

Nell'applicazione dei modelli linguistici N-gram risulta complicato stimare le probabilità N-gram di $P(\varphi_k|\varphi_{k-N+1}, \dots, \varphi_{k-1})$. La risoluzione più intuitiva è data dalla stima della massima verosimiglianza (Bishop, 2006), con cui si cerca la configurazione delle probabilità non conosciute per massimizzare la funzione di verosimiglianza su un insieme di dati che vengono utilizzati per addestrare il sistema. Nel caso generale della stima della massima verosimiglianza per i modelli linguistici N-gram, si stima la seguente probabilità condizionata

$$P(\varphi_k|\varphi_{k-N+1}, \dots, \varphi_{k-1}) = \frac{C(\varphi_{k-N+1}, \dots, \varphi_{k-1}, \varphi_k)}{C(\varphi_{k-N+1}, \dots, \varphi_{k-1})} \quad (3)$$

dove $C(\varphi_{k-N+1}, \dots, \varphi_{k-1})$ è la frequenza della sequenza di parole $\varphi_{k-N+1}, \dots, \varphi_{k-1}$ nel testo di addestramento.

4.1 Probabilistic topic models

I modelli tematici costituiscono una categoria di modelli generativi utilizzati per scoprire la struttura semantica sottostante di una raccolta di documenti. L'idea nasce da Deerwester *et al.* nel 1990, i quali concepirono l'indicizzazione semantica latente

(LSI), in cui viene eseguita la decomposizione del singolo valore per scoprire strutture statistiche attraverso e nei documenti in uno spazio dimensionale inferiore. Per interpretare gli argomenti latenti scoperti, nel 1999 Hoffman pensò ai pLSI (probabilistic latent semantic indexing), in cui un documento è modellato come una miscela di argomenti latenti e ogni argomento è modellato come una distribuzione multinomiale sulle parole. Tuttavia, i modelli pLSI non sono modelli generativi completi e per questo è stato introdotto nel 2003 da Blei *et al.* un modello probabilistico bayesiano completo, chiamato allocazione Dirichlet latente (LDA), grazie al quale la proporzione dell'argomento in ogni documento si presume sia estratta dalla distribuzione di Dirichlet nello stesso corpus. Questo modello è un traguardo importante che ha aperto molte possibilità allo sviluppo ulteriore di vari modelli tematici generativi. È servito da trampolino di lancio per molti altri modelli tematici per diversi tipi di dati testuali (Steyvers *et al.*, 2004; Lafferty, 2007; Hong and Davison, 2010; Lin and He, 2009; Wang and Blei, 2011; Zhao *et al.*, 2011; Jo and Oh, 2011; Wang *et al.*, 2011).

4.2 *L'analisi semantica latente probabilistica*

L'analisi semantica latente probabilistica (pLSA), nota anche come indicizzazione semantica latente probabilistica (pLSI) è un modello generativo il cui obiettivo è l'analisi delle co-occorrenze di dati per scoprire la struttura semantica sottostante dei dati. Per descrivere formalmente il pLSI, è opportuno prima introdurre alcune notazioni e terminologie.

Una parola φ è l'unità di base definita in un vocabolario a dimensione fissa, indicizzato da 1 a V . Un documento è una sequenza di parole di lunghezza N , indicata con $d = (\varphi_1, \varphi_2, \dots, \varphi_N)$. Un corpus è una collezione di M documenti, indicati con $D = (d_1, d_2, \dots, d_M)$. Nella pLSI si presume che un corpus contenga un insieme di k argomenti latenti, ciascuno dei quali è modellato come una distribuzione multinomiale sul vocabolario (ad esempio $p(\varphi|\beta_i)$), dove β_i è il parametro di distribuzione dell'argomento i .

Ogni parola in un documento è generata da un singolo argomento indicizzato da z e parole diverse in un documento possono essere generate da diversi argomenti. Un'ipotesi importante è che, una volta assegnati gli argomenti $z = (z_1, z_2, \dots, z_N)$ alle parole del documento d , le parole sono indipendenti dall'indice del documento. Di conseguenza, la probabilità congiunta del documento d e le sue parole $\varphi_1, \varphi_2, \dots, \varphi_N$ possono essere calcolate come

$$P(d, \varphi_1, \varphi_2, \dots, \varphi_N) = P(d) \prod_{i=1}^N \sum_{z_i} P(\varphi_i | z_i) P(z_i | d). \quad (4)$$

La scomposizione della probabilità congiunta di un documento e delle sue parole nella pLSI può essere descritta dal seguente processo generativo:

1. Per ogni $d \in D$, campionare d da $d \sim p(d)$.
2. Per generare ogni parola $\varphi_i \in d$,
 - a. campionare l'assegnazione di argomenti z_i da $z_i \sim p(z|d)$;
 - b. campionare parola φ_i da $\varphi_i \sim p(\varphi|\beta, z_i)$.

L'indicizzazione semantica latente probabilistica attenua l'assunto per cui l'intero corpus contiene un solo argomento e ogni parola nei documenti viene campionata da quell'argomento, introducendo k argomenti latenti in una data collezione e consente ad ogni documento di essere una miscela di questi k argomenti. In sostanza, ogni documento è rappresentato come un elenco di proporzioni di miscelazione per queste componenti miscelate, come $p(z|d)$ e quindi ridotto ad una distribuzione di probabilità su un insieme fissato di argomenti. Quelle proporzioni di miscelazione possono essere considerate come una rappresentazione di dimensione inferiore di un documento, che può anche essere considerato un'utile conoscenza riguardo alla copertura degli argomenti in ogni documento.

Il modello pLSI è stato utilizzato come punto di partenza da molti ricercatori. Brants *et al.* (2002) lo hanno adoperato per eseguire la segmentazione di documenti basati su argomenti. Mei *et al.* (2007, 2005) se ne sono serviti per modellare le sfaccettature e le opinioni nei blog in rete e scoprire schemi di temi evolutivi dal testo.

Il modello pLSI ha due parametri da stimare, che sono la distribuzione della parola di un argomento dato i , $p(\varphi|\beta_i)$ e le proporzioni degli argomenti nel documento dato d , $p(z|d)$. A causa dell'esistenza di variabili latenti, la stima della massima verosimiglianza non è più applicabile e per stimare questi parametri viene utilizzato l'algoritmo di massimizzazione delle aspettative (algoritmo EM). Quest'ultimo è un metodo iterativo per trovare la massima probabilità o il massimo di una stima dei parametri nei modelli che dipendono da variabili latenti non osservate.

Il numero dei parametri nel modello pLSI cresce linearmente con la dimensione del corpus e non è chiaro come assegnare la probabilità ad un documento al di fuori dell'addestramento. Per queste limitazioni è stato introdotto il modello di allocazione Dirichlet latente (LDA) per imporre un presupposto generativo completo sul processo di generazione del documento.

4.3 Allocazione Dirichlet latente

Il modello di allocazione Dirichlet latente (LDA), proposto da Blei *et al.* nel 2003 introduce una distribuzione Dirichlet condivisa sulle proporzioni degli argomenti in ciascun documento per controllare il numero di parametri in un modello

argomentativo. La proposta argomentativa $p(z|\theta, d)$ nel documento d è modellata come distribuzione multinomiale parametrizzata da un vettore k -dimensionale θ che si presume essere estratto da una distribuzione di Dirichlet, in cui α è il parametro di concentrazione and $\Gamma(\cdot)$ è la funzione Gamma:

$$p(\theta|\alpha) = \frac{\Gamma(\sum_{i=1}^k \alpha_i)}{\prod_{i=1}^k \Gamma(\alpha_i)} \prod_{l=1}^K \theta_l^{\alpha_l - 1}. \quad (5)$$

Il processo generativo del documento specificato dal modello LDA può essere descritto come segue:

1. Per ogni $d \in D$, campionare θ da $\theta \sim \text{Dir}(\alpha)$.
2. Per ogni parola $\varphi_i \in d$,
 - a. campionare l'assegnazione di argomenti z_i da $z_i \sim p(z|\theta, d)$;
 - b. campionare parola φ_i da $\varphi_i \sim p(\varphi|\beta, z_i)$.

La corrispondente probabilità congiunta di parole φ , le assegnazioni di argomenti latenti z e la proporzione dell'argomento latente nel documento specificato da un modello LDA possono essere calcolate come:

$$p(\varphi, z, \theta|\alpha, \beta) = p(\theta|\alpha) \prod_{n=1}^N p(\varphi_n|\beta, z_n) p(z_n|\theta). \quad (6)$$

Il modello LDA si basa su un'ipotesi Bayesiana gerarchica a due livelli nel processo di generazione del documento: la proposta argomentativa θ viene estratta da una distribuzione di Dirichlet e l'assegnazione specifica dell'argomento di ogni parola viene estratta da una distribuzione multinomiale specificata da θ . Ciò fornisce un vantaggio computazionale aggiuntivo, che facilita l'inferenza a posteriori. Rispetto al modello pLSI, la proposta argomentativa θ viene modellata come una variabile latente, piuttosto che come un parametro del modello, in modo tale da avere il numero dei parametri del modello LDA indipendente dal corpus di addestramento.

Anche questo modello ha posto le basi ad una serie di applicazioni, come ad esempio, le dinamiche temporali di una distribuzione di parole sugli argomenti in un dato corpus (Blei and Lafferty, 2006), la supervisione continua (Mcauliffe and Blei, 2008) e la supervisione discreta (Zhu *et al.*, 2009; Ramage *et al.*, 2009). Nel 2006 Teh *et al.* hanno introdotto un altro strato di gerarchia Bayesiana sulla generazione del parametro di Dirichlet, in modo tale che la proprietà di raggruppamento dei documenti possa essere acquisita. A causa dell'accoppiamento tra la variabile continua θ e la variabile discreta z in un documento, l'inferenza a posteriori nel modello LDA risulta essere più vantaggiosa rispetto a quella del modello pLSI. I due metodi inferenziali maggiormente usati sono il campionamento di Gibbs (Griffiths and Steyvers, 2004) e

l'inferenza variazionale (Blei *et al.*, 2003) ed entrambi sfruttano la coniugazione tra la distribuzione di Dirichlet e la distribuzione multinomiale per facilitare il calcolo.

Una implementazione parallela al modello LDA per la raccolta di documenti su larga scala è stata studiata da Smola e Narayanamurthy (2010), da Zhai *et al.* (2012) e da Wang *et al.* (2009).

5. Modelli upstream e downstream

Nel 2008 Mimno e McCallum hanno classificato i modelli generativi per la sentiment analysis in modelli upstream e modelli downstream, in base alla loro particolare assunzione di dipendenza tra le etichette dei sentimenti s , l'assegnazione degli argomenti z e la parola osservata φ in un dato documento.

I modelli upstream presumono che, per generare una parola $\varphi_{d,n}$ in un dato documento d , per prima cosa bisogna decidere la polarità del sentiment $s_{d,n}$ di questa parola. Tale polarità determina l'assegnazione dell'argomento $z_{d,n}$ per questa parola. Generalmente, questi modelli identificano il sentiment con etichette discrete e presumono che ci siano diverse proporzioni di argomento sotto diverse etichette di sentiment.

I modelli downstream, invece, presumono che l'etichetta del sentiment $s_{d,n}$ sia determinata dall'assegnazione dell'argomento $z_{d,n}$ parallelamente alla parola $\varphi_{d,n}$. In sostanza questi modelli risultano maggiormente flessibili nella costruzione del sentiment.

Il modo in cui viene specificata la dipendenza rappresenta appunto la diversità tra questi due modelli. Intuitivamente, nei modelli upstream argomenti e parole sono potenzialmente dipendenti dalla variabile del sentiment e in questo modo può essere considerata "a monte" la sua influenza su altre variabili catturate direttamente dal modello. Nel modello downstream si presume che la variabile del sentiment dipenda dagli argomenti e il modello tenta di acquisire come altre variabili influenzino la variabile del sentiment, considerata, appunto, "a valle". Trattando il sentiment come variabile risposta della variabile argomento, ci sono diversi modi per costruire il sentiment, anche nel caso delle classificazioni numeriche, che sarebbero difficili da costruire con un modello upstream.

5.1 Modelli upstream per la sentiment analysis

I modelli upstream presumono che, per generare una parola in un documento di testo, sia necessario per prima cosa campionare un'etichetta di sentiment latente per poi

campionare una etichetta di argomento rispetto a questa categoria di sentiment e, infine, campionare la parola da questo argomento scelto.

Tra i modelli generativi upstream per la sentiment analysis è molto utilizzato il modello Topic-Sentiment Mixture (TSM), proposto da Mei *et al.* nel 2007, che si basa sul modello pLSI. Oltre a presupporre che un corpus sia composto da k argomenti con sentiment neutro, il TSM introduce due modelli di sentiment aggiuntivi, uno per le opinioni negative ed uno per le opinioni positive. Nel TSM si suppone che i modelli di sentiment siano ortogonali ai modelli argomentativi, nel senso che assegnerebbero alte probabilità alle parole generiche che sono spesso usate per esprimere polarità di sentiment, mentre i modelli argomentativi assegnerebbero alte probabilità alle parole che rappresentano i contenuti degli argomenti con opinioni neutrali.

Nel TSM viene introdotto un nuovo concetto chiamato “tema” che presenta tre componenti: parole neutre, parole positive e parole negative in ogni documento. La combinazione tra modelli argomentativi e modelli di sentiment crea un tema riguardo ad un particolare aspetto con una certa polarità di sentiment in un dato documento. Questa combinazione varia tra i diversi documenti per riflettere distinte polarità di sentiment degli utenti verso lo stesso aspetto. Una volta determinati i temi, un documento viene modellato come una miscela di temi e il resto del processo di generazione segue quello del modello pLSI. Poiché il modello TSM si basa sul modello pLSI, soffre delle sue limitazioni, come l’overfitting, che possono, però, essere arginate con le ipotesi del modello LDA.

Nel 2009 Lin e He hanno proposto il modello Joint Sentiment and Topic (JST), nel quale si presume che un corpus contenga $S \times k$ argomenti, dove S è il numero delle categorie di sentiment, positive, negative e neutrali. La combinazione di argomenti e sentiment viene modellata come un prodotto cartesiano tra modelli di argomento e modelli di sentiment, in una maniera piuttosto simile all’interpolazione lineare presunta nel modello TSM. Il modello JST prima campiona una etichetta di sentiment e poi campiona l’assegnazione dell’argomento e la parola dalle distribuzioni corrispondenti. Per generare un documento con il modello JST, è necessario prima campionare una miscela di sentimenti per questo documento da una distribuzione condivisa di Dirichlet e in ciascuna categoria di sentiment campionare un argomento mescolando una parte da un’altra distribuzione di Dirichlet. La proporzione dell’argomento in ogni documento viene modellata come vettori S -dimensionali, che consentono diverse miscele di argomenti in diverse categorie di sentiment.

Anche il modello Aspect and Sentiment Unification (ASUM) proposto da Jo e Oh nel 2011 utilizza lo stesso presupposto generativo del modello JST, ma, per rafforzare la coerenza dell’argomento e del sentiment all’interno del documento, presume che

tutte le parole in una frase condividano lo stesso argomento e assegnazione del sentiment. Anche nel modello ASUM viene applicata la stessa inferenza a posteriori del modello JST, che considera una frase come unità di base dell'inferenza.

Una diversa variante del modello generativo upstream per l'analisi del sentiment è stata proposta nel 2010 da Zhao *et al.* Il principio della Massima Entropia (ME) è stato introdotto nel modello LDA per controllare la selezione delle parole da argomenti generici, argomenti specifici e argomenti suddivisi in base alle opinioni. Nel modello ME-LDA una determinata parola può essere generata da diverse categorie di argomenti. L'assegnazione di una parola particolare ad una di queste categorie viene controllata dal modello ME basato su caratteristiche discriminative estratte dal contenuto della parola.

La maggior parte dei modelli generativi upstream più diffusi utilizzati per l'analisi del sentiment trattano l'etichetta del sentiment come una variabile latente e il sentiment, definito "prior" viene impiegato per iniettare la polarità del sentiment nei modelli. Sebbene un tale approccio offra flessibilità nell'identificare opinioni distinte su singole parole, si necessita di una grande conoscenza del sentiment per garantire risultati di analisi soddisfacenti.

5.2 Modelli downstream per la sentiment analysis

I modelli downstream invertono il presupposto di generazione tra le etichette di sentiment e l'assegnazione di argomenti latenti: per generare un documento di testo è necessario selezionare prima le assegnazioni degli argomenti in questo documento e campionare le parole e le etichette di sentiment rispetto a questi argomenti. Un esempio di modello generativo downstream per l'analisi del sentiment è il modello supervisionato LDA (sLDA) di McAuliffe e Blei (2008). Il presunto processo di generazione del contenuto di testo nel modello sLDA è identico a quello assunto nel modello LDA. Oltre alla generazione del documento, il modello sLDA presuppone che la variabile risposta y sia ricavata da una distribuzione Gaussiana con media $\eta^T \bar{z}$ e deviazione standard σ , dove $\bar{z} = \frac{1}{N} \sum_{n=1}^N z_n$, come ad esempio il vettore medio delle assegnazioni di argomenti nel documento d . Con questo presupposto il modello sLDA può essere utilizzato come modello di regressione per costruire le classificazioni di opinioni nei documenti di testo. Si può generare la variabile risposta y con un modello lineare generalizzato, come ad esempio un modello logistico, per costruire classi discrete di sentiment. Come nel modello LDA, l'inferenza variazionale può essere applicata nel modello sLDA per l'inferenza a posteriori.

Nel 2009 Zhu *et al.* hanno introdotto l'idea di addestrare il margine massimo nei modelli sLDA per una migliore prestazione predittiva.

Nel 2010 Boyd-Graber e Resnik hanno generalizzato il modello sLDA per ottenere una sentiment analysis olistica tra le lingue. Nel loro modello gli argomenti sono organizzati secondo una struttura semantica condivisa e l'etichetta del sentiment in un dato documento viene costruita come una variabile risposta rispetto all'assegnazione dell'argomento. Di conseguenza il modello identifica simultaneamente come i concetti multilingue vengano raggruppati in argomenti tematicamente coerenti e come gli argomenti associati al testo siano collegati alla classificazione del sentiment.

Nel 2011 Wang e Blei hanno esteso la sLDA ad un ambiente congiunto, dove il filtraggio basato sulle classificazioni delle opinioni degli utenti possa essere ottenuto in uno spazio di argomenti latenti.

Lin e He nel 2012 hanno eseguito un'interessante riparametrizzazione del modello JST per trasformare il loro modello originale JST upstream in un modello downstream, denominato Reverse-JST. In questo modello si presume che per generare la sequenza di parole in un dato documento è necessario prima campionare l'assegnazione di argomenti per poi campionare la categoria di sentiment rispetto all'argomento selezionato e, infine, selezionare una parola da questa combinazione di sentiment.

Un'importante filone di ricerca nei modelli generativi downstream per la sentiment analysis si concentra sulla comprensione delle opinioni che riguardano gli aspetti, che possono essere classificati come polarità di sentiment degli utenti sugli argomenti latenti in un dato documento. Titov e McDonald nel 2008 hanno sviluppato un modello generativo basato sul modello LDA, chiamato Multi-Aspect Sentiment (MAS) per la modellazione congiunta del contenuto del testo e la classificazione degli aspetti per il riepilogo dei sentiment. Nella loro soluzione, due tipi di argomenti, locali e globali, sono modellati esplicitamente e si presume che ogni frazione all'interno di un documento sia una miscela di argomenti locali e globali. In base alle assegnazioni degli argomenti latenti, si presume che le classificazioni delle espressioni vengano determinate da un modello di regressione logistica, che considera le assegnazioni degli argomenti e la sequenza di parole come input. Rispetto al modello sLDA, MAS consente la comprensione del sentiment con una previsione dettagliata delle opinioni. E' molto importante che le etichette del sentiment siano note nella fase di addestramento, altrimenti la loro assenza può costituire una grande limitazione del modello.

6. Conclusioni

La generalità e la robustezza dei modelli generativi della sentiment analysis sono qualità che li rendono di utile applicazione a tutti i dati testuali e strumenti importantissimi per la costruzione di sistemi intelligenti per la comprensione delle opinioni. In questo contributo si è voluto passare in rassegna questi importanti modelli per sottolineare la loro utilità nell'inferenza di variabili latenti su opinioni dettagliate e la loro efficacia nelle analisi congiunte di dati testuali e valutazioni numeriche associate.

Riferimenti bibliografici

- Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*, XX, 738 p. ISBN: 978-0-387-31073-2. Springer.
- Blei, D.M.; Ng, A.Y.; Jordan, M.I. (2003). Latent Dirichlet allocation. *The Journal of Machine Learning Research* 3: 993–1022.
- Blei, D.M.; Lafferty, J.D. (2006). Dynamic topic models. In *Proceedings of the 23rd International Conference on Machine Learning*, 113–120.
- Blei, D.M., Lafferty, J.D. (2007). A correlated topic model of science. *The Annals of Applied Statistics* 1(1): 17–35.
- Bolasco, S.; Bisceglia, B.; Baiocchi, F. (2004). Estrazione automatica d'informazione dai testi, *Mondo Digitale* 1: 27- 43.
- Boyd-Graber, J.; Resnik, P. (2010). Holistic sentiment analysis across languages: Multilingual supervised latent Dirichlet allocation. In *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing (EMNLP '10)*, 45–55, Stroudsburg. Association for Computational Linguistics.
- Brants, T.; Chen, F.; Tsochantaridis, I. (2002). Topic-based document segmentation with probabilistic latent semantic analysis. In *Proceedings of the Eleventh International Conference on Information and Knowledge Management*, 211–218.
- Deerwester, S. C.; Dumais, S.T.; Landauer, T. K.; Furnas, G. W.; Harshman, R. A. (1990). Indexing by latent semantic analysis, *JASIS* 41(6): 391-407.
- Griffiths, T.L.; Steyvers, M. (2004). Finding scientific topics. *Proceedings of the National Academy of Sciences*, 101(suppl 1): 5228–5235.

- Hiemstra, D.; Kraaij, W. (1998). Twenty-one at TREC7: ad-hoc and cross-language track. In *Proceedings of The Seventh Text REtrieval Conference (TREC 1998)*, Gaithersburg, 174–185, 9–11 Nov 1998.
- Hofmann, T. (1999). Probabilistic latent semantic indexing. In *Proceedings of the 22nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, 50–57.
- Jo, Y.; Oh, A.H. (2011). Aspect and sentiment unification model for online review analysis. In *Proceedings of the Fourth ACM International Conference on Web Search and Data Mining*, 815–824.
- Katz, S.M. (1987). Estimation of probabilities from sparse data for the language model component of a speech recognizer. *IEEE Transactions on Acoustics, Speech and Signal Processing* 35(3): 400–401.
- Lin, C.; He, Y.; Everson, R. ; Rüger, S. (2012). Weakly supervised joint sentiment-topic detection from text. *IEEE Transactions on Knowledge and Data Engineering* 24(6): 1134–1145.
- McAuley, J.; Leskovec, J. (2013). Hidden factors and hidden topics: Understanding rating dimensions with review text. In *Proceedings of the 7th ACM Conference on Recommender Systems*, 165–172.
- Mcauliffe, J.D.; Blei, D.M. (2008). Supervised topic models. In *Advances in Neural Information Processing Systems*, 121–128.
- Mei, Q.; Zhai, C. (2005). Discovering evolutionary theme patterns from text: An exploration of temporal text mining. In *Proceedings of the Eleventh ACM SIGKDD International Conference on Knowledge Discovery in Data Mining*, 198–207.
- Mei, Q.; Ling, X.; Wondra, M.; Su, H.; Zhai, C. (2007). Topic sentiment mixture: Modeling facets and opinions in weblogs. In *Proceedings of the 16th International Conference on World Wide Web*, 171–180.
- Mimno, D.; McCallum, A. (2008). Topic models conditioned on arbitrary features with Dirichlet-multinomial regression. *The 24th Conference on Uncertainty in Artificial Intelligence*, 411–418.
- Moghaddam, S.; Ester, M. (2011). ILDA: Interdependent LDA model for learning latent aspects and their ratings from online product reviews. In *Proceedings of the 34th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 665–674.
- Ponte, J.M.; Croft, W.B. (1998). A language modeling approach to information retrieval. In *Proceedings of the 21st Annual International ACM SIGIR Conference on (and Development in Information Retrieval (SIGIR '98)*, 24–28 Aug 1998, Melbourne, 275–281.

- Rabiner, Lawrence R.; Biing-Hwang Juang. (1993). *Fundamentals of speech recognition*. Upper Saddle River: Prentice-Hall.
- Russel, S.; Norving, P. (2005). *Intelligenza artificiale. Un approccio moderno*, Prentice Hall 2 – Pearson Education Italia.
- Smola, A.; Narayanamurthy, S. 2010. An architecture for parallel topic models. *Proceedings of the VLDB Endowment* 3(1–2): 703–710.
- Steyvers, M.; Smyth, P.; Rosen-Zvi, M.; Griffiths, T. (2004). Probabilistic author-topic models for information discovery. In *Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 306–315.
- Teh, Y.W.; Jordan, M.I.; Beal, M. J.; Blei, D. M. (2006). Hierarchical Dirichlet processes. *Journal of the American Statistical Association* 101(476).
- Titov, I.; McDonald, R.T. (2008b). A joint model of text and aspect ratings for sentiment summarization. In *ACL*, vol. 8, 308–316. Citeseer
- Wang, C.; Blei, D.M. (2011). Collaborative topic modeling for recommending scientific articles. In *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 448–456.
- Wang, H.; Lu, Y.; Zhai, C. (2010). Latent aspect rating analysis on review text data: A rating regression approach. In *Proceedings of the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 783–792.
- Wang, H.; Lu, Y.; Zhai, C. (2011). Latent aspect rating analysis without aspect keyword supervision. In *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 618–626.
- Wang, Y.; Bai, H.; Stanton, M.; Chen, W.Y.; Chang, E.Y. (2009). PLDA: Parallel latent Dirichlet allocation for large-scale applications. In *Algorithmic Aspects in Information and Management*, 301–314. Springer.
- Zhai, C.; Lafferty, J. (2001a). Model-based feedback in the language modeling approach to information retrieval. In *Proceedings of the Tenth International Conference on Information and Knowledge Management*, 403–410.
- Zhai, K.; Boyd-Graber, J.; Asadi, N.; Alkhouja, M.L. (2012). Mr. LDA: A flexible large scale topic modeling package using variational inference in mapreduce. In *Proceedings of the 21st International Conference on World Wide Web*, 879–888.
- Zhao, W.X.; Jiang, J.; Yan, H.; Li, X. (2010). Jointly modeling aspects and opinions with a maxent-lda hybrid. In *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing (EMNLP '10)*, Stroudsburg, 56–65. Association for Computational Linguistics.

- Zhao, W.X.; Jiang, J.; Weng, J. He, Lim, E.P.; Yan, H.; Li, X. (2011). Comparing twitter and traditional media using topic models. In *Advances in Information Retrieval*, 338–349. Springer.
- Zhu, J.; Ahmed, A.; Xing, E.P. (2009). Medlda: Maximum margin supervised topic models for regression and classification. In *Proceedings of the 26th Annual International Conference on Machine Learning*, 1257–1264.

Comparison of different multivariate calibrations and ensemble methods for estimating selected soil properties with vis-NIR reflectance spectroscopy

Davide Fragnito^{1*}, Natalia Leone², Valeria Ancona²,
Domenico Vitale³, Antonio Lucadamo⁴

¹Master's Graduate in Statistical, Actuarial and Financial Sciences,

²Water Research Institute, National Research Council, Viale Francesco de Blasio, 5, Bari,

³CMCC Foundation, Euro-Mediterranean Center on Climatic Change, Viterbo,

⁴DEMM, Department of Law, Economics, Management and quantitative Methods,
University of Sannio, Via delle Puglie, 53, Benevento

Abstract: Sustainable soil management requires a correct assessment of soil chemical and physical properties. Historically, this has been gained through conventional laboratory analyses, which are considered costly and time-consuming, particularly when a large number of soil samples need to be analysed. An alternative, faster and less expensive, approach is based on the use of reflectance spectroscopy in the vis-NIR domain. This approach implies the calibration of predictive models that relate the spectral reflectance to soil properties. The goodness of the models can be particularly influenced by the multivariate methods used. In this article, we compare the performance of different multivariate and statistical ensemble methods for estimating some basic soil properties, such as sand, silt, clay, and organic carbon in the specific pedo-environmental conditions of an important agricultural area in southern Italy.

Keywords: vis-NIR reflectance spectroscopy, prediction of soil properties, multivariate and statistical ensemble methods.

1. Introduction

Soil is one of the main natural resources. It contributes to basic human needs like food, clean water, and clean air, and is a major carrier for biodiversity (Keesstra et al., 2016). From here, the need to preserve this resource (soil) to ensure sustainable

* Corresponding author: davide.fragnito.1994@gmail.com.

and shared prosperity to humanity (FAO, 2015) through sustainable agricultural and non-agricultural uses. Sustainable soil management cannot disregard a correct assessment of its chemical and physical properties and their variability in space and time.

Historically our understanding of the soil system and assessment of its properties has been gained through conventional laboratory analysis (Viscarra-Rossel et al., 2006). The latter, although usefully and practically irreplaceable for detailed investigations, are costly and time-consuming, thus not very suitable when large numbers of soil samples need to be analysed, as, for example, in large soil surveys, or for high-resolution soil mapping and precision agriculture. Hence, the need to develop alternative techniques for soil analyses.

In recent years, vis-NIR reflectance spectroscopy has been shown to be a useful technique for the measurement of various soil properties (Lucadamo and Leone, 2015; Lucadamo et al., 2020). Compared to conventional analytical methods, vis-NIR spectroscopy is faster, cheaper, and non-destructive; it requires less sample preparation, with less or no chemical reagents, is highly adaptable to automated and in situ measurements, and has the potential to analyse various soil properties simultaneously (Viscarra-Rossel et al., 2006; McCarty et al., 2002; Vasques et al., 2008). Reflectance spectroscopy refers to the measure of spectral reflectance (Milton, 1987), i.e., the ratio of the electromagnetic radiation reflected by a soil surface to that which impinges on it (Drury, 1993). Since the characteristics of the radiation reflected from a material are a function of the material's properties, observations of soil reflectance can provide information on the properties and state of the soil (Irons et al., 1989). The reflectance spectra of soil in the vis-NIR are largely non-specific due to the overlapping absorption of soil constituents. This characteristic lack of specificity is compounded by scatter effects, caused by soil structure or specific constituents, such as quartz. All of these factors result in complex absorption patterns that need to be mathematically extracted from the spectra and correlated with soil properties. Therefore, multivariate statistics are required to mathematically extract complex absorption patterns and to correlate these patterns with the measured soil properties for calibration (Martens and Næs, 1989; Stenberg et al., 2010; Araújo et al., 2014; Xu et al., 2018). The selection of the multivariate statistic methods, along with that of proper instrumentation, accessories and optical probe design (Mouazen et al., 2009), improved spectra filtering and pre-processing (Maleki et al., 2008), are essential factors for successful calibration of predictive models (Mouazen et al., 2010; Nawar et al., 2016).

A large number of multivariate calibration methods have been used to relate vis-NIR reflectance spectra with measured soil properties (e.g., Viscarra-Rossel et al., 2006; Janik et al., 2009; Mouazen et al., 2010; Stevens et al., 2010; Viscarra-Rossel and Behrens, 2010; Vohland et al., 2011; Shi et al., 2015; Araújo et al., 2014; Kuang et al., 2015; Were et al., 2015). However, none of these proposed calibration techniques have achieved universal acceptance because a calibration model that works well for one application may be unacceptable for another (Xu, 2018). The specificity of the pedo-environment, besides the choice of the pre-processing methods, may also influence the selection of the statistical calibration methods, being a soil a complex and heterogeneous system.

This study aims to explore the performances of different multivariate and statistical ensemble methods for estimating some basic soil properties, such as sand, silt, clay, and organic carbon (OC) contents, within the specific pedo-environmental conditions of an important, irrigated area of southern Italy. Namely, the compared statistical methods are: Partial Least Squares Regression (PLS), Regression Tree (RT), Bagging and Random Forest algorithm (B, RF), Boosting Regression (BR), Artificial Neural Network (ANN), Multivariate Adaptive Regression Splines (MARS). The remaining part of the article is organized as follows: in section 2, all the statistical methods used for the analysis are introduced; section 3 describes data collection and material; in section 4 the results are synthesized; some concluding remarks are shown in section 5.

2. Some theoretical aspects

2.1 Partial Least Squares Regression (PLSR)

Partial Least Squares Regression is by far the most used multivariate statistical method in the field of vis-NIR reflectance spectroscopy (Gholizadeh et al., 2016; Leone et al., 2012; Leone et al., 2019; Vibhute et al., 2018; Viscarra-Rossel et al., 2006; Cozzolino and Moron, 2003; Wang et al., 2013; Volkan Bilgili et al., 2010; Lee et al., 2009; Viscarra-Rossel and Behrens, 2010; Kuang et al., 2015; Wetterlind et al., 2008; Viscarra-Rossel and Lark, 2009; Brown et al., 2006; Stevens et al., 2013; Dunn et al., 2002; Fystro, 2002; Mouazen et al., 2007). This method was proposed by H. Wold for the modeling of data sets in terms of chains of matrices (path models), suggesting a procedure named NIPALS (Non-linear Iterative Partial Least Squares) to estimate the parameters (Wold, 1973). Later, other groups led by S. Wold and H. Martens popularized the use of this method for chemical applications by

slightly modifying the PLS model with only two matrices containing the explanatory variables (X) and the response variables (Y) to deal with complicated data sets where ordinary regression was difficult or impossible to apply. Several authors (Wold et al., 1993; Wold et al., 2004) started to interpret PLS as the Projection to Latent Structures, providing a more descriptive meaning.

There are two basic approaches, named PLSR1 and PLSR2. In PLSR1, one calibration model is considered for y or separate calibration models are built for each column in Y . With PLS2, one calibration model is built for all columns of Y simultaneously. PLSR was first proposed for analysing NIR spectra by Wold et al. (Wold et al., 1983), who derived an algorithm with orthogonal scores. Successively, Martens (Martens, 1985) and Martens and Naes (Martens and Naes, 1989) proposed a PLSR algorithm with orthogonal loadings. Moreover, Helland (Helland, 1998) showed the equivalence between these proposals for the PLSR1 algorithms, while the geometry of PLS has been explored in depth by Phatak and de Jong (Phatak and De Jong, 1997). For the purpose of this paper, we refer only to the PLSR1 algorithm.

PLSR finds the linear (or polynomial) relationships between a centred response variable vector y and a matrix of centred predictors X expressed as $y=f(X)+E$. PLS regression seeks then to provide a statistical model based on the reduction of the space spanned by the often-large number of correlated predictors in a lower-dimensional space generated by derived PLS components. These components reflect the information in the X -variables that are of relevance for modelling and predicting the response variable y . The link is then obtained by the following decompositions that lead to orthogonal scores and non-orthogonal loading vectors (Wold's algorithm):

$$\mathbf{X} = \mathbf{t}_1\mathbf{p}'_1 + \mathbf{t}_2\mathbf{p}'_2 + \dots + \mathbf{t}_K\mathbf{p}'_K + \mathbf{E}_K = \mathbf{TP}' + \mathbf{E}_K$$

$$\mathbf{y} = \mathbf{t}_1\mathbf{q}_1 + \mathbf{t}_2\mathbf{q}_2 + \dots + \mathbf{t}_K\mathbf{q}_K + \mathbf{f}_K = \mathbf{Tq} + \mathbf{f}_K$$

where t is a vector of scores calculated by $t_k = X_{k-1}w_k$ with scaled weights w_k and $T = [t_1, \dots, t_k]$, p are the spectral loadings, q the chemical loadings and E and f are the predictor and response variable residuals, respectively, of the estimated effect for the k -th factor ($k=1, \dots, K$). Wold et al. use, for achieving these solutions, the well-known non-linear iterative partial least squares (NIPALS) algorithm for centred X and y data: let $X_0 = X$ and $Y_0 = Y$, the orthogonal scores $\{t_1, \dots, t_k\}$ are then iteratively obtained, where the basic k -th step of the algorithm is given by:

1. Compute the scaled weight vector $w_k = cX'_{k-1}Y_{k-1}/Y'_{k-1}Y_{k-1}$ with c scaling factor;
2. Compute the orthogonal score $t_k = X_{k-1}w_k$;
3. Compute the residuals $X_k = (I - P_{t_k})X_{k-1} = (I - P_{t_k})X$ and

$y_k = (I - P_{t_k})y_{k-1} = (I - P_{t_k})y$ where $P_{t_k} = t_k(t'_k t_k)^{-1}t'_k$ and $P_{T_k} = T_k(T'_k T_k)^{-1}T'_k$ are the orthogonal projection operators onto t_k and the subspace spanned by $\{t_1, \dots, t_k\}$, respectively.

The number of factors to use in PLSR model may be determined through leave-one-out cross-validation. The optimal number of factors should allow the modelling of as much as possible of the correlation between X and y without overfitting y . Then, for the selected number of factors, one calculates the final linear regression coefficients, $b = W(W'P)^{-1}q$ (where W is the weight matrix) and $b_0 = \bar{y} - \bar{x}'b$ to be used in the predictor $\hat{y}_i = b_0 + x_i b$ where x_i is the new spectrum. The well-known Martens algorithm is instead based on the factorization:

$$\begin{aligned} X &= \tilde{t}_1 \tilde{w}'_1 + \dots + \tilde{t}_K \tilde{w}'_K + E_K = \tilde{T} \tilde{W}' + E_K \\ y &= \tilde{t}_1 \tilde{q}_1 + \dots + \tilde{t}_K \tilde{q}_K + f_K = \tilde{T} \tilde{q} + f_K \end{aligned}$$

which uses a non-orthogonal score matrix \tilde{T} , i.e. $\tilde{T}' \tilde{T}$ is a non-diagonal matrix, with orthogonal loadings and where the scores $\{t_1, \dots, t_k\}$ are iteratively obtained. The basic k -th step of this algorithm (with $X_0 = X$ and $Y_0 = Y$) is given by:

- 1) Compute the weight vector $\tilde{w}_k = \tilde{X}'_{k-1} \tilde{Y}'_{k-1}$;
- 2) Compute the non-orthogonal score $\tilde{t}_k = \tilde{X}'_{k-1} \tilde{w}'_k / \tilde{w}'_k \tilde{w}_k$ and set $\tilde{T}_k = [\tilde{t}_1, \dots, \tilde{t}_k]$;
- 3) Compute the regression coefficients \tilde{q}_k of y in \tilde{T}_k given by $\tilde{q}_k = (\tilde{T}'_k \tilde{T}_k)^{-1} \tilde{T}'_k \tilde{Y}_k$;
- 4) Compute the residuals $\tilde{X}_k = \tilde{X}_{k-1} - \tilde{t}_k \tilde{w}'_k$ and $\tilde{y}_k = y - \sum_{j=1}^k \tilde{q}_j \tilde{t}_j$;

such to obtain the above X decomposition $X = \tilde{T} \tilde{W}' + \tilde{E}_k$, and where the score vectors $\tilde{t}_k = X \tilde{w}_k / \tilde{w}'_k \tilde{w}_k$ and t_k span the same vectorial space. The regression coefficient vector is finally given as a simple least square solution $\tilde{b} = \tilde{W}(\tilde{W}' X' X \tilde{W})^{-1} \tilde{W}' X' y$ providing the same coefficients as the previous PLS1 formula.

We remark that the latter algorithm, giving the non-orthogonal score vectors, does not provide the problems that Pell (Pell et al., 2007) has recently highlighted for the NIPALS results about their possible inconsistency with respect to model spaces for residual-based outlier detection and prediction purpose. See Ergon (Ergon, 2009) for a re-interpretation of the NIPALS results, which solves the PLSR inconsistency problem.

We highlight that in this paper all computations have been performed by using the software for the chemometric analysis of spectroscopic data called "ParLeS"

(Viscarra-Rossel, 2008). This software implements the most used form of the PLSR1 algorithm, which produces orthogonal scores, and provides several statistical tools to assist the researcher in performing and interpreting the analysis results. For example, the number of samples (rows) to leave out in the cross-validation may be any integer selected by the user and the accuracy of the cross-validation is given by the root-mean-square error (RMSE). Moreover, the goodness of fit is given by R^2 and Q^2 statistics, which give the upper and lower bounds, of how the model well explains the data and predicts new observations. For the selection of an optimal parsimonious PLSR model (i.e., one that represents the variability in the data without causing it to overfit) the Akaike Information Criterion (AIC) (Akaike, 1973) is also provided by ParLeS where N is the sample size and m is the number of model parameters, in this case, the number of factors. A sorted VIP (Variable Importance for Projection) data table, and the percent variation in each of the x and y -data that is explained by each of the PLSR factors, are also given, where the VIP index is computed as:

$$\text{VIP}_j(\mathbf{k}) = K \sum_k w_{jk}^2 \left(\text{SSY}_k / \text{SSY}_{\text{tot}} \right)$$

where $\text{VIP}_j(\mathbf{k})$ is the importance of the j -th predictor variable based on a model with k factors, w_{jk} is the corresponding loading weight of the j -th variable in the k -th PLSR factor, SSY_k is the explained sum of squares of y by a PLSR performed with the only k -th factor, SSY_{tot} is the total sum of squares of y , and K is the total number of predictor variables. The reader is directed to Viscarra-Rossel (Viscarra-Rossel, 2008) for a full description of ParLeS and the algorithms it implements.

2.2 Regression Trees (RT)

An alternative algorithm for analysing the relationship between variables is the "Classification And Regression Trees" (CART). Even if the basic idea is the same, in this paper it has been preferred to separate the classification tree treatment from the tree regression, also by virtue of the fact that in the dataset in our possession, the variable y is a continuous random variable.

The regression tree constructs an H tree from the root node h_1 , by performing a succession of splits, or divisions, of the full set of observations, to make the units more homogeneous in terms of response variable y . The algorithm used to build the tree follow an approach of step-by-step optimization. To understand how it works, we must break down the deviance as follows:

$$D = \sum_{i=1}^n [y_i - \hat{f}(x_i)]^2 = \sum_{h=1}^J \{ \sum_{i \in R_h} (y_i - \hat{c}_h)^2 \} = \sum_h D_h \quad (1)$$

where $\hat{f}(x_i)$ is the predicted value for response variable y ; c_h is the arithmetic mean of the observed y_i having component x_i falling in the subinterval; J is the global number of nodes and R_1, \dots, R_J are rectangles in the p -dimensional sense. The growth process of the tree starts with the root node $h1$ (so $J=1$; $R_j = \mathbb{R}^p$; $D = \sum_{i=1}^n [y_i - M(y)]^2$ with $M(\cdot)$ average operator). And proceeds iteratively according to the following scheme:

- Once a rectangle R_h is chose, the appropriate value of c_h is the average of the corresponding values $\hat{c}_h = M(y_i: x_i \in R_h)$;
- If we subdivide the region into two parts the deviance is replaced by $D_h^* = \sum_{i \in R_{h'}} (y_i - \hat{c}_h')^2 + \sum_{i \in R_{h''}} (y_i - \hat{c}_h'')^2$ with a gain of $g_h = D_h - D_h^*$.
- We can inspect all p explanatory variables and, for each of them, all the possible points of subdivision, selecting the variable and its point of subdivision that maximize g_h .

The algorithm stops when all the leaves contain a number of sample elements that is less than a preassigned value, or when the relative fall of deviance is less than a prefixed threshold. A large tree is obviously not useful, so the branches of little importance have to be pruned. For this reason, a cost-complexity function can be considered:

$$C_\alpha(J) = \sum_{h=1}^J D_h + \alpha J$$

where α is a non-negative penalty parameter. For each α there is a unique smallest tree minimizing $C_\alpha(J)$. The algorithm sequentially eliminates one leaf at a time. At each step the leaf for which elimination causes the smallest increase in $\sum_h D_h$ is selected. The question is to choosing α : generally the cross-validation is used. Trees are frequently used in practice, but it is important to underline their advantages and disadvantages. They have in fact a logical simplicity and are easy to communicate; the step function has a simple, compact mathematical formulation in terms of information to be stored; there is a speed of computation and the possibility to use discrete and categorical variables; a robust forms of deviance can be used; not particularly complicated variations can be introduced, which allow for missing values, in both tree construction and prediction; the method automatically selects the important variables. On the other hand, there is instability of results and difficulty in upgrading the tree; difficulty of approximating some mathematically simple function; procedures of statistical inference are not available and it is not simple to evaluate the order of importance of variables remaining in the pruned tree. (Breiman et al., 1984; Ripley, 1996; Venables & Ripley, 1997; AA.VV., 1995).

2.3 Bootstrap Aggregating (B)

Regression trees suffer from high variance, which means that if you were to divide the calibration dataset into two random parts and fit a regression tree to both halves, the results that could be achieved would be quite different. Conversely, a procedure with low variance will produce similar results when applied repeatedly to separate datasets; linear regression tends to have a low variance if the ratio of n (number of observations) to “ p ” (number of predictors) is moderately large.

Bootstrap AGGREGATING, or BAGGING, is a general procedure for reducing the variance of a statistical learning method and is particularly useful, and often used, in the context of regression trees, although applied in different works (Gholizadeh et al., 2016; Viscarra-Rossel and Behrens, 2010). Briefly remembering that the variance of the mean of observations \bar{Z} of a set of n independent observations Z_1, \dots, Z_n , each with variance σ^2 , is equal to the ratio of variance to the number of observations, it is statistically valid to say that the average of a series of observations reduces variance.

Hence a natural way to reduce variance and simultaneously increase the accuracy of predictions of a statistical learning method is to take many sub-samples of the "training" of the model from the population, build a separate forecasting model using each train-set, set and calculate the average of the resulting predictions. In other words, you could calculate $\hat{f}^1(x), \hat{f}^2(x), \dots, \hat{f}^B(x)$ using separate B training sets, mediate them into $\hat{f}_{\text{avg}}(x)$ obtaining a single low-variance statistical learning model, expressed by:

$$\hat{f}_{\text{avg}}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}^b(x)$$

In general, this method is not widely applied, considering that everyone does not have easy access to multiple training datasets, for various reasons, such as the impossibility of replicating that phenomenon or for purely economic issues or lack of time available. In contrast, the bootstrap technique obtains reproductions of samples from the individual training dataset, generating different B "bootstrapped" train samples. Finally, the statistical method is trained through the bootstrap set b -th in such a way as to obtain $\hat{f}^{*b}(x)$, and mediate all the predictions to obtain $\hat{f}_{\text{bag}}(x)$ like:

$$\hat{f}_{\text{bag}}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}^{*b}(x)$$

So, this is the bagging procedure and to apply this algorithm to regression trees, you simply generate B trees using B bootstrapped train samples and mediate the resulting predictions. These trees are allowed to grow and are not subject to

"Pruning" procedures, i.e. pruning final observations. It has also been shown in other works (James et al., 2013) that to observe significant improvements in accuracy, the bagging procedure requires hundreds of tree replications in a single procedure.

2.4 Random Forest (RF)

Random forest is a method also applied in Viscarra-Rossel and Antoine Stevens (Viscarra-Rossel and Behrens, 2010; Stevens et al., 2013) and provides an improvement over trees developed with the bagging procedure through a small modification that decorrelates the trees. As in bagging, we build a series of decision trees on bootstrapped training samples. But when you build these decision trees, whenever a division in a tree is considered, a random sample of m predictors is chosen as candidates apart from the complete set of predictors. A new sample of m predictors is taken at each division, generally $m \approx \sqrt{p}$, which means the number of predictors considered in each division is roughly equal to the square root of the total number of predictors. In other words, in the construction of a random forest, with each division of the tree, the algorithm is not even allowed to consider most of the available predictors.

This may sound crazy, but it has intelligent logic. Suppose you have a very powerful predictor in your dataset, along with a number of other moderately strong predictors. So, in the collection of bagging trees, most or all trees will use this strong predictor in the upper-division (James et al., 2013). As a result, all bagging trees will look quite similar to each other, resulting in highly correlated predictions. Unfortunately, the average of many highly correlated quantities does not lead to a large reduction in variance as the average of many unrelated quantities. In particular, this means that the bagging algorithm will not result in a substantial reduction in variance on a single tree in this setting. Random forests overcome this problem by forcing each division to consider only a subset of the predictors.

Therefore, on average $\frac{p-m}{p}$ divisions will not even consider the strong predictor, and therefore other predictors will have a better chance. You can think of this process as a decoration of the trees, thus making their average less variable and, therefore, more reliable. The main difference between bagging and random forests is the choice of the size of the predictor subset of m size. For example, if a random forest is constructed using m plus p , this is simply the same as bagging. On the data used, the random forests they use lead to a reduction in error compared to the $m \approx \sqrt{p}$ bagging procedure.

Using a small value of m in building a random forest will generally be useful when we have a large number of related predictors. As with bagging, random forests

do not adapt too much to the increase in the number of B iterations, but in practice, a value large enough to allow the error rate to stabilize has stabilized (James et al., 2013).

2.5 Boosting Regression (BR)

Like the bagging algorithm, boosting, or "enhancement," is a general approach that can be applied to many statistical learning methods for regression or classification. Here we limit our discussion on incentive to the context of regression trees, as approached before by Gholizadeh, Brown and Stevens (Gholizadeh et al., 2016; Brown et al., 2006; Stevens et al., 2013).

Remember that bagging involves creating multiple copies of the original training dataset using bootstrap, adapting a separate decision tree to each copy, and then combining all the trees to create a single predictive model. In particular, each tree is based on a bootstrap dataset, independent of other trees. Boosting works in a similar way, except that trees are grown sequentially: each tree is grown using information from previously grown trees. The upgrade does not involve bootstrap sampling, but each tree adapts to a modified version of the original dataset.

Considering the approach of the regressive technique, such as bagging, boosting also involves the combination of a large number of decision-making trees, $\hat{f}^1(x), \hat{f}^2(x), \dots, \hat{f}^B(x)$. The algorithm that governs the regressive approach of enhancement could thus be summarized as follow:

1. Consider the relation between the dependent variable y and the explicative ones: $y = f(x)$.
2. Set $\hat{f}(x) = 0$ and $r_i = y_i$ for each " i " in the train dataset; r_i and y_i are the residuals and the value of the response variable for the generic observation i , respectively.
3. For $b = 1, 2, \dots, B$ repeat the following sub-process:
 - a) Fit a tree \hat{f}^b with d splits (i.e. $d-1$ terminal nodes) to the dataset train.
 - b) Update \hat{f} to add in a reduced version of the new regression tree: $\hat{f}(x) \leftarrow \hat{f}(x) + \lambda \hat{f}^b(x)$ where λ is a shrinking parameter.
 - c) Update residuals as follows: $r_i \leftarrow r_i - \lambda \hat{f}^b(x)$.
4. Get the boosted model, as:

$$\hat{\mathbf{f}}(\mathbf{x}) = \sum_{b=1}^B \lambda \hat{\mathbf{f}}^b(\mathbf{x})$$

The idea behind this procedure is this: unlike adapting a single large decision tree to the dataset, which equates to forced and potentially excessive data fitting (it could

be affected by "overfit", that is, by overfitting the data by calibrating excessively correctly and validating less than enough), the enhancement approach allows for slow learning. Given the current model, it was preferred to adopt a regression tree to the remnants of the model. That is, we adapt a tree using the current residuals, instead of the variable Y , as an answer. You then add this new decision tree in the adapted function to update the residues. Each of these trees can be quite small, with a few terminal nodes, determined by the d parameter in the algorithm. By adapting small trees to the residues, we notice an improvement, albeit slow, of \hat{f} areas where it does not work well. The shrinking parameter λ further slows down the process, allowing more and different-shaped trees to attack the residues. In general, slow-learning statistical learning approaches tend to work well. Note that in upgrading, the construction of each tree depends heavily on the trees that have already been cultivated. It can then be summarized that boosting has three optimization parameters:

- 1) The number of B trees. Unlike bagging and random forests, boosting can be oversized to the data if B is too large, although this oversizing tends to occur slowly if at all.
- 2) The shrinking parameter λ , a small positive number, which controls the speed with which it learns boosting. Typical values are 0.01 or 0.001 and the right choice may be the problem. A very small value of λ requires the use of a very large value of B to achieve good performance.
- 3) The number of divisions in each tree, which controls the complexity of the boosting set. Often d plus 1 works well, in case each tree is a "stump", consisting of a single division. In this case, the boosting set adapts to an additive model, because each term involves only a single variable. More generally, the d parameter can be interpreted as the depth of interaction and controls the interaction order of the boosting model, as d divisions can involve, at most, d variables. This highlights a difference between enhancement and random forests: in boosting, given that the growth of a particular tree considers others that have already been trained, it can be trusted to trust the condition that smaller trees are sufficiently adequate even in interpretation. For example, the use of stumps, mentioned above, leads to an additive pattern (James et al., 2013).

2.6 Artificial Neural Network (ANN)

The Artificial Neural Network (ANN) represents an artificial reproduction of a biological neural network of a human brain, including "neurons", nerve cells that are interconnected in a real network, and applied to predict soil contents, above all, by Kuang (Kuang et al., 2015).

However, it is important to note that in such a widespread network not all interconnections have the same specific weight in terms of importance; in fact, some have a high priority, associated with greater weight, than others. Like biological networks, artificial neural networks also have interconnected neurons and a pattern that faithfully reports the structure of a biological neural network.

Historically, the earliest ANNs are The Perceptron, proposed by Rosenblatt (Rosenblatt, 1958) and the Artron due to R. Lee (Lee, 1959). Then the Adaline (Adaptive Linear Neuron) and The Madaline (Many Adaline), due to Widrow et al. (Widrow et al., 1960, 1988). The first one is an artificial neuron also known as the ALC (adaptive linear combiner), the ALC being its principal component. The second one is an ANN (network) formulation based on the Adaline above, but it is a multilayer NN. Principles of the above four neurons are common building blocks in almost all ANN architectures.

Four major multi-layer general-purpose network architectures are:

- The Back-Propagation network: a multi-layer Perceptron-based ANN, giving an elegant solution to hidden-layers learning (Rumelhart et al., 1986). Its computational elegance stems from its mathematical foundation that may be considered as a gradient version of Richard Bellman's Dynamic Programming theory (Bellman, 1954)
- The Hopfield Network (Hopfield, 1982): this network is different from the earlier ANNs in many important aspects, especially in its recurrent feature of employing feedback between neurons. Hence, although several of its principles have been incorporated in ANNs based on the earlier four ANNs, it is to a great extent an ANN-class in itself. Its weight adjustment mechanism is based on the AM principle
- The Counter-Propagation Network (Hecht-Nielsen, 1987): Kohonen's Self-Organizing Mapping (SOM) is employed to facilitate unsupervised learning, utilizing the WTA principle to economize computation and structure.
- The LAMSTAR (LargeMemory Storage And Retrieval) network: a Hebbian network that uses a multitude of Kohonen SOM layers and their WTA principle. It is unique in its employs these by using Kantian-based Link-Weights (Graupe and Lynn, 1969) to link different layers (types of stored information) The link weights allow the network to simultaneously integrate inputs of various dimensions or nature of representation and incorporating correlation between input words. Furthermore, the network incorporates (graduated) forgetting in its learning structure and it can continue running uninterrupted when partial data is missing.

In this paper we use the Back-propagation method.

2.7 *Multivariate Adaptive Regression Splines (MARS)*

Many of the classic regression models have only linear aspects, but they can be adapted to non-linear models in the data by manually adding nonlinear terms to the model; however, in order to do so, the analyst must know a priori the specific nature of non-linearities and interactions. Alternatively, there are many inherently non-linear algorithms. When using these models, the exact shape of the non-linearity should not be explicitly known or specified before the model is formed. Rather, these algorithms will look for and discover non-linearities and interactions in data that help maximize predictive accuracy.

An example of such algorithms is the Multivariate Adaptive Regression Spline (MARS) (Friedman, 1991), an algorithm that automatically creates a linear pattern that sometimes provides an intuitive approach to the non-linearity after grasping the concept of multiple linear regression. They provide a cost-effective approach to capturing non-linear relationships in your data by evaluating breakpoints (nodes) similar to step functions. The procedure evaluates each data point for each predictor as a node and creates a linear regression model.

The MARS procedure will first search for the single point through the x-value range where two different linear relationships between Y and X reach the smallest error. For a single node, the hinge function is of the type:

$$y = \begin{cases} \beta_0 + \beta_1(a_1 - x) & x < a_1 \\ \beta_0 + \beta_1(x - a_1) & x > a_1 \end{cases}$$

Once the first node is found, the search continues for a second node. This procedure continues until many nodes are found, producing (potentially) a highly non-linear forecast equation. Including many nodes can allow you to adapt a really good relationship with the available training data, but it could lead you not to generalize, and therefore predict, very well with new and/or unknown data. Therefore, once you have identified the complete set of nodes, you can sequentially remove nodes that do not contribute significantly to predictive precision. This process is known as "pruning" and has been used to find the optimal number of nodes. There are two important optimization parameters associated with the MARS model: the maximum degree of interactions and the number of terms maintained in the final model. A grid search is necessary to identify the optimal mix of hyperparameters, i.e., the different combinations of interaction complexity and the number of terms to keep in the final model.

The advantages of MARS models are numerous. First, they naturally manage mixed types of predictors (quantitative and qualitative), consider all possible binary

partitions of categories for a quantitative predictor in two groups, thus generating a pair of indicative functions for the two categories. These templates also require minimal functionality design and automatically select the features. For example, because they scan each predictor to identify a subdivision that improves predictive accuracy, non-informational features will not be selected. What is more, highly correlated predictors do not prevent predictive accuracy as much as OLS models.

However, one drawback of MARS models is that they are generally slower to train. Because the algorithm analyses each predictor value for potential breakpoints, computational performance can be affected by both increases in the number of observations and the number of variables. In addition, although related predictors do not necessarily hinder the model's performance, they can make it difficult to interpret. When two features are "almost perfectly" related, the algorithm will essentially select the first one that occurs when scanning features. Therefore, choosing one at random, the related function probably will not be included because it does not add any explanatory power to the analysis. (Gene et al., 1979). In soil environment, MARS are used to predict textures and contents by many authors (Volkan Bilgili et al., 2010; Viscarra-Rossel and Behrens, 2010; Nawar et al., 2016; Stevens et al, 2013).

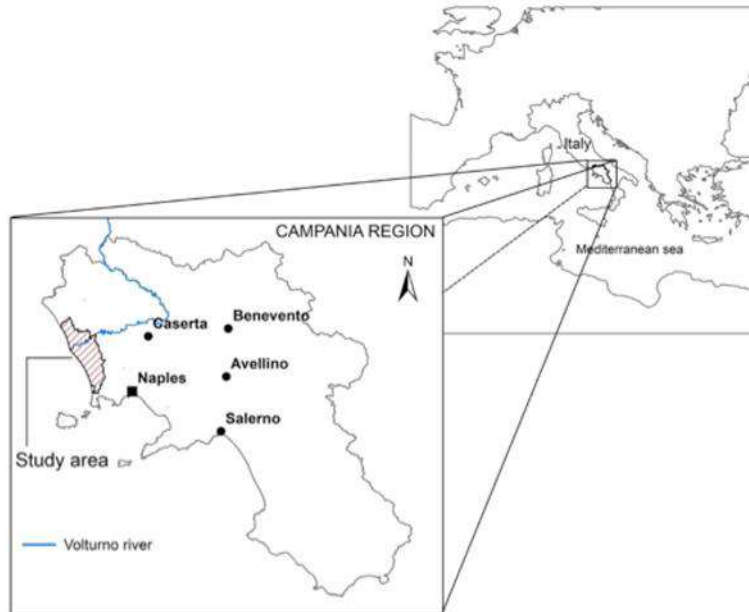
3. Materials and methods

3.1 Study area and soil sampling

The area under investigation (Figure 1) is located in the north-western part of the Campania Region, in southern Italy (Coord. 41°01'00'' N, 13°58'00'' E), within a fertile agricultural land, mainly devoted to irrigated vegetal crops and fruit trees (Geoportale Regione Campania, 2019). The climate is typically Mediterranean, with the wettest period between late autumn (October–November) and early spring (March–April). Temperature and potential evapotranspiration temperatures show an inverse trend compared to rainfall, with the highest values during summer (June–August).

The dominant soils types are Gleyic, Gleyic-Vertic, Calcari-Gleyic and Calcari-Fluvic Cambisols, and Calcaric Gleysols (Di Gennaro, 2002). For this study, an existing soil database, made available CNR-ISAFoM, was used. The database contains information on soil organic carbon (OC) and particle size distribution (sand, silt, and clay), used in our application. Information about OC was available for ninety-six samples, while those for sand, silt, and clay contents were available only for eighty-two samples.

Figure 1. Localisation of the study area in southern Italy.



The analytical data regarded surface soil samples randomly collected in 1999 within the study area, air-dried, and ground to a size fraction passing a 2 mm sieve. Soil organic carbon and texture were determined according to the Italian Official Methods for Soil Analysis (MIPAF, 2000). Namely, total clay (soil separate with < 0.002 mm particle diameter) and silt (soil separate with 0.002 to 0.05 mm particle diameter) contents were determined with the pipet method. Total sand content (soil separate with 0.05 to 2.0 mm particle diameter) was determined by wet sieving; OC content was determined using Walkley-Black methods.

3.2 Vis-NIR spectroscopy

The diffuse vis-NIR spectral reflectance was measured in the laboratory, on a residual fraction of soil samples, under controlled light conditions, using the procedure described in Leone et al. (2019). Noisy portions of the measured reflectance spectra, between 350 and 399 nm and between 2451 and 2500 nm, were removed, leaving spectra in the range of 400-2450 nm for the analysis. The resulting reflectance spectra were normalised, using the continuum removal approach (Clark and Roush, 1984). To this end, a convex hull was fitted over the original spectral curve, then the absorption spectrum was calculated by taking the ratio between the original reflectance spectrum and the enveloping curve (Van der Meer, 1999; De Jong, 1992).

3.3 Statistical calibrations

The selected statistical calibrations were performed to predict the investigated soil properties from reflectance spectra, using both the software “ParLeS” (Viscarra-Rossel, 2008) and “R x64 3.6.3” (R Core Team, 2020). All models performed to calibrate the spectral data with the reference (laboratory) soil data have employed two-thirds of the available samples for calibration and the remaining third for independently validating them. For each variable, the selection of samples was carried out as follows: first, the samples were sorted following ascending order of the variable, then, sequentially, every two samples were taken for calibration and the third for validation.

To enhance the predictive power of these statistical calibration models, spectroscopic data were transformed and pre-processed prior to data analysis, with the aim of removing undesired variation in the data (Eriksson et al., 2006). In this study, we assessed all the transformation and pre-processing methods, either alone or in combination, before calibrations.

The combination of the following procedures provided the best results: reflectance (R) to absorbance (A) transformation ($A = \log 1/R$), wavelet detrending, median filtering, second derivative of absorbance, and data enhancement (mean centre). In particular, reflectance to absorbance transformation reduces nonlinearities (Viscarra-Rossel, 2008), while wavelet detrending (Daubechies, 1992) corrects light scattering variation and baseline. The median filter, in addition (Viscarra-Rossel, 2008), reduces the effects of random spectral noise, thereby providing smoother spectra. Lastly, the second derivative removes additive and linear baseline effects (Burger and Geladi, 2007), while amplifying absorption features, which are indicative of the contents of the soil materials. Mean centring is a commonly used method of data enhancement to reduce redundant information and better evaluate differences.

Leave-one-out cross-validation (Efron and Tibshir, 1994) was then used to determine the number of factors to retain in the calibration models. To select the optimal cross-validated calibration model, we computed the root mean square error (RMSE) of predictions:

$$\text{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_{\text{pred}} - y_{\text{ob}})^2}$$

in which N is the sample population size, y_{pred} is the predicted value, and y_{ob} is the observed value. In this case, the model with the lowest RMSE is selected. However, a more parsimonious model, i.e., a model with fewer factors representing the variability

in the data set, without causing overfitting, is preferred. For that purpose, the optimal selection of factors can be based on the penalizing Akaike Information Criterion (AIC) (Akaike, 1969; Li et al., 2002):

$$\text{AIC} = N \log (\text{RMSE}) + 2m.$$

in which N is the sample population size and m is the number of model parameters (i.e., the number of factors). This criterion is applied after have verified that residuals have zero mean normal distribution.

To evaluate the accuracy of models, the adjusted coefficient of determination (R_{adj}^2) and the relative percent deviation (RPD), i.e., the ratio of the standard deviation of analysed data (i.e., the soil properties) to RMSE, was performed. In accordance with previous studies (Williams, 1987; Viscarra-Rossel, 2007) the quality of predictions expressed by RPD was classified as follows: $\text{RPD} < 1.0$ indicates very poor model/predictions and their use is not recommended; RPD between 1.0 and 1.4 indicates poor model/predictions where only high and low values are distinguishable; RPD between 1.4 and 1.8 indicates fair model/ predictions which may be used for assessment and correlation; RPD values between 1.8 and 2.0 indicates good model/ predictions where quantitative predictions are possible; RPD between 2.0 and 2.5 indicates very good, quantitative model/ predictions, and $\text{RPD} > 2.5$ indicates excellent model/predictions. RPD statistic is also carried out to assess the performance of validation using the independent data set.

4. Results and discussion

4.1 Descriptive statistics of soil properties

The investigated soil variables were statistically described in terms of minimum, maximum, mean, coefficient of variation (CV), and skewness. Furthermore, a log-transformation was performed for those variables that did not follow a normal distribution. Summary statistics of calibration and validation subsets are reported in Table 1. Organic carbon content ranges from 2.71 to 215.6 g Kg^{-1} , and is on average moderate (21.5 g Kg^{-1}). A slight difference in the average values can be observed between the calibration (22.5 g Kg^{-1}) and validation (19.6 g Kg^{-1}) sub-sets. Skewness always exhibits high values: 4.57 g Kg^{-1} for the whole dataset; 4.52 and 3.74 g Kg^{-1} , for the calibration and validation sub-sets, respectively, thus indicating a significant deviation from the normal distribution.

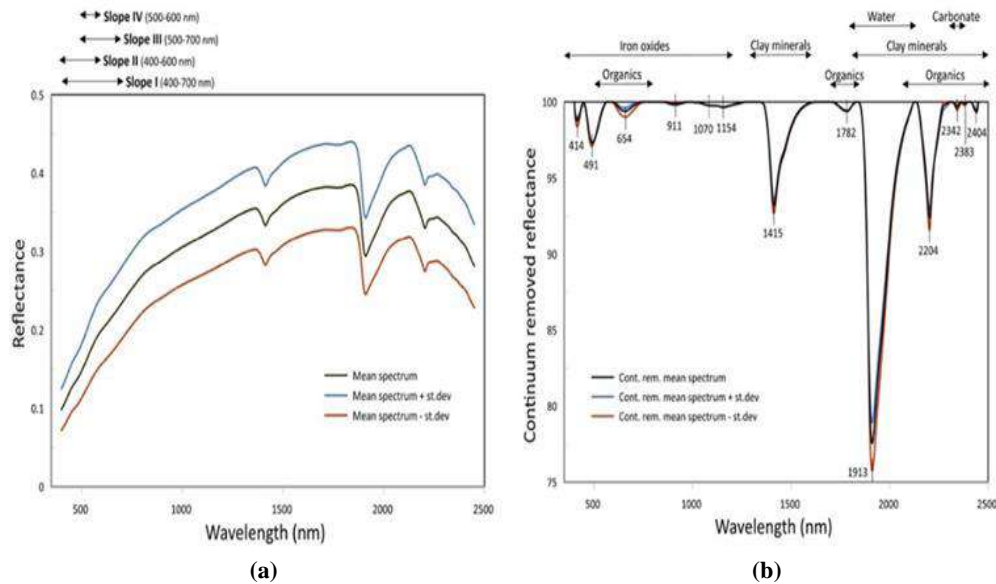
Table 1. Descriptive statistics of the selected soil properties for calibration and validation datasets.

		n	Mean	Range	CV	Skewness
OC (g Kg ⁻¹)	Calibration	64	22.5	4.4 – 215.6	1.43	4.52
	Validation	32	19.6	2.8 – 121.4	1.12	3.74
Sand (g Kg ⁻¹)	Calibration	54	419.8	80.0 – 940.0	0.53	0.59
	Validation	28	423.6	70.0 – 950.0	0.56	0.64
Silt (g Kg ⁻¹)	Calibration	54	201.1	10.0 – 370.0	0.39	-0.35
	Validation	28	201.1	10.0 – 390.0	0.43	-0.23
Clay (g Kg ⁻¹)	Calibration	54	377.4	50.0 – 730.0	0.46	-0.14
	Validation	28	378.9	10.0 – 770.0	0.50	-0.09

Soil separates, i.e., the size groups of mineral particles, is dominated by the sand, (421.1 g Kg⁻¹), on average, followed by clay (377.7 g Kg⁻¹) and silt (201.1 g Kg⁻¹) fractions. The dominant, basic soil textural classes are: clay, clay-loam, sandy-clay-loam, and sandy-loam. Extreme and mean values for all, sand, silt, and clay calibration and validation subsets are similar. Skewness was consistently low, thus indicating, for these variables, a frequency distribution close to the normal distribution. Differences between calibration and validation sub-sets are minimal, and the CV is moderate for both these variables. Skewness is consistently low. Considering that the mean and coefficient of variation (CV) for the calibration and validation sets are comparable for all the considered soil properties, the selection of both datasets can be considered representative (Ding et al., 2018). Figure 2 shows the average soil spectrum and the relative continuum removed reflectance of the investigated soil samples and their standard deviation.

The average spectrum (Figure 2a) shows a typical convex shape and a moderate overall reflectance. The dispersion of the spectral intensity, as measured by the standard deviation, was evident. Many studies demonstrated that different soil properties, especially particle size distribution and organic carbon content, may affect the overall reflectance (Stenberg et al., 2010). Changes in slopes of different ranges in the visible region are also observed. Various studies have related visible reflectance slope to soil organic matter content (Summers et al., 2011). The average continuum removed spectrum (Figure 2b) shows several absorption bands across the entire vis-NIR region, which can be related to clay minerals, organic matter, iron oxides, water, and carbonate contents (Stenberg et al., 2010; Leone A.P., 2000).

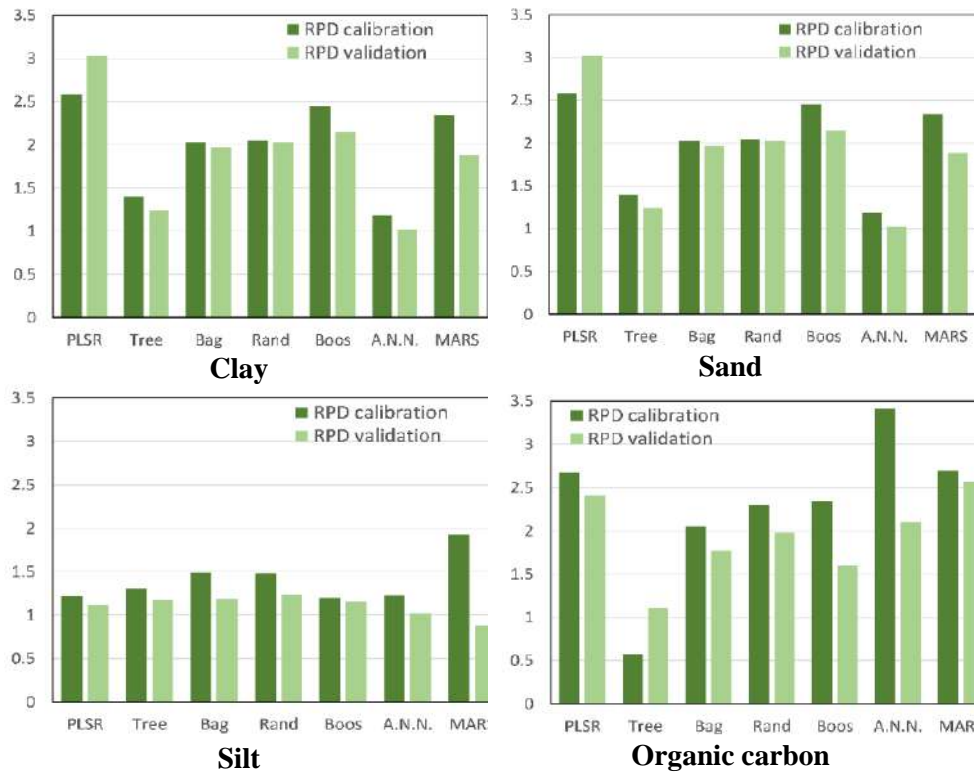
Figure 2. Mean of spectral reflectance **(a)** and continuum removed spectral reflectance **(b)** of sampled soils. In **(a)** the position of spectral ranges where the various visible reflectance slopes were calculated; in **(b)** the approximate positions of some fundamental soil constituents are shown.



4.2 Multivariate and ensemble calibrations

The capability of vis-NIR reflectance spectroscopy to predict the investigated soil properties among different statistical methods is summarised in Table 2, and shows that, as a general rule and considering both calibration and validation results, PLS gives the best results. The comparison among different models is immediately evident in Figure 3. However, the response of those methods in both Table 2 and Figure 3 gives slightly different results depending on the variable considered. In any case, clay and organic carbon are the best-predicted variables. Specifically, PLSR applied to two-thirds of the available sample set revealed good correlations between soil reflectance spectra and the considered soil properties, except for silt content. Based on the RPD values, the calibration models were excellent for log-OC (RPD = 2.67) and clay (RPD = 2.59) and good for sand (RPD = 1.95). For clay, models including 5 factors, based on the RMSE and AIC values, allowed to attain a cross-validation between predicted and measured data with R_{adj}^2 of 0.855, 0.845 and 0.731, for log-OC, clay and sand, respectively. For silt, the only calibration possible performed a poor model, with an R_{adj}^2 of 0.314 and RPD of 1.22.

Figure 3. Histograms of calibration and validation RPD of all statistical models for soil variables.



In some cases, increasing the number of factors gave slightly higher coefficients of regression (R^2), but increased RMSE, thus reducing the stability of the calibration models (i.e., leading to over-fitting) (Vågen et al., 2006; Wise et al., 2003). Therefore, we selected the most parsimonious model in terms of number of factors, based on the values of AIC and RMSE.

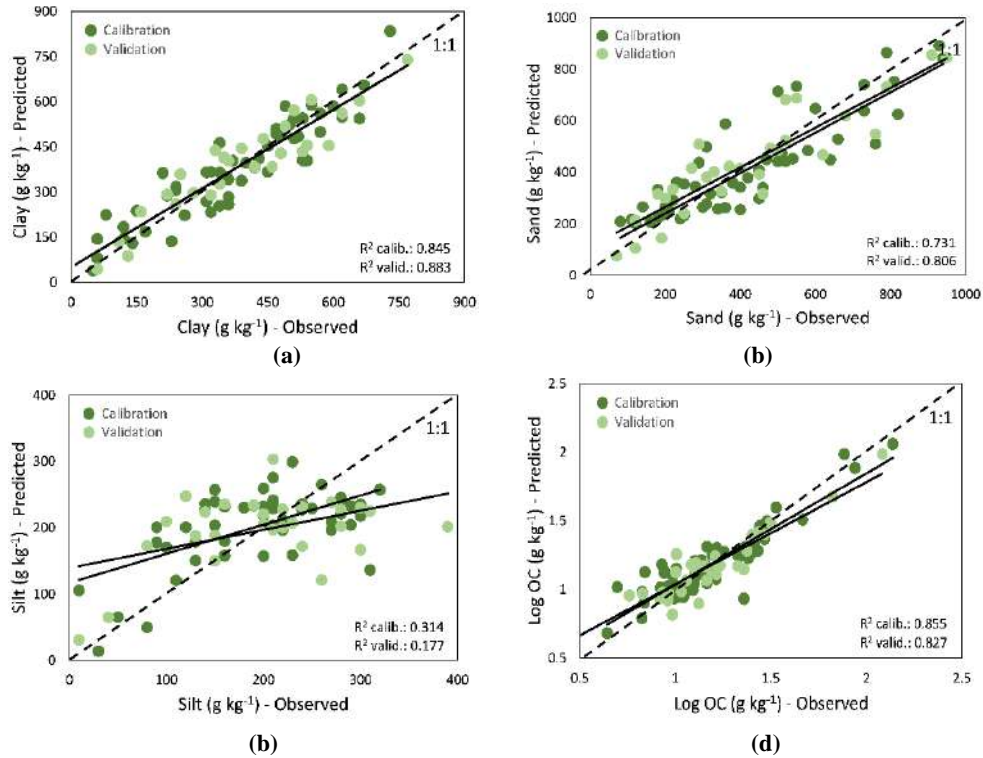
Leave-one-out calibration models constructed through vis-NIR reflectance spectroscopy and PLSR are empirical; therefore, validations of these models are better performed using a data set that is independent of the one used for calibration (Volkan Bilgili et al., 2010). Validation using the remaining one third of the available samples indicated excellent models for the prediction of clay ($R_{adj}^2 = 0.883$ and RPD = 3.03), very good models for log-OC ($R_{adj}^2 = 0.827$ and RPD = 2.41) and for sand ($R_{adj}^2 = 0.806$ and RPD = 2.10), and a poor model for silt ($R_{adj}^2 = 0.177$ and RPD = 1.11).

Table 2. RMSE and RPD for different methods applied on calibration and validation sample for Clay, Sand, Silt and OC contents in soil samples.

Variables	Models	Calibration sample		Validation sample	
		RMSE	RPD	RMSE	RPD
Clay	PLSR	67.61	2.59	62.46	3.03
	Tree	125.18	1.40	152.51	1.24
	Bagging	86.31	2.03	95.76	1.97
	Random Forest	85.47	2.05	93.07	2.03
	Boosting	71.28	2.45	88.05	2.15
	A.N.N.	146.72	1.19	185.56	1.02
	M.A.R.S	74.70	2.34	100.54	1.88
Sand	PLSR	114.61	1.95	102.40	2.32
	Tree	120.02	1.86	134.17	1.77
	Bagging	116.83	1.91	117.16	2.03
	Random Forest	120.48	1.86	113.18	2.10
	Boosting	118.17	1.89	121.67	1.96
	A.N.N.	186.15	1.20	233.66	1.02
	M.A.R.S	117.48	1.90	135.04	1.76
Silt	PLSR	64.99	1.22	77.19	1.12
	Tree	60.20	1.31	73.19	1.18
	Bagging	52.99	1.49	73.03	1.19
	Random Forest	53.45	1.48	69.60	1.24
	Boosting	65.80	1.20	74.53	1.16
	A.N.N.	64.44	1.23	84.99	1.02
	M.A.R.S	41.09	1.93	97.90	0.88
OC	PLSR	0.12	2.67	0.13	2.41
	Tree	0.20	0.58	0.28	1.11
	Bagging	0.15	2.05	0.17	1.77
	Random Forest	0.14	2.30	0.15	1.98
	Boosting	0.13	2.35	0.19	1.61
	A.N.N.	0.09	3.42	0.15	2.10
	M.A.R.S	0.12	2.70	0.12	2.57

To make the reader more comfortable with the results of PLSR prediction, scatterplots of the predicted vs measured values for these properties are shown in Figure 4. In this plot, the values of the R_{adj}^2 and regression's straight lines of both the datasets are also highlighted. In order to make exhaustive the discussion about the results of this work, the outcomes of the other models cannot be overlooked.

Figure 4. Scatterplots of observed vs predicted soil properties for calibration and validation data sets in PLSR.



In particular, considering singularly all different models examined through the statistical-computational environment R (R Core Team, 2020), there are some models with behaviour similar to PLSR. One of these models is MARS, that has performed very good/excellent values in terms of RPD in Clay and OC ($2.3 \leq RPD_{cal} \leq 2.7$) and good ones in Sand and Silt. It is remarkable that MARS is better than PLSR to predict OC, as evidenced by the excellent values of RMSE and RPD in both calibration and validation datasets. In the validation phase, this model returns good outputs ($1.8 \leq RPD_{val} \leq 2.6$), excluding Silt, in which predictions are not recommended ($RPD_{val} \leq 1$). The second model in order of goodness of predictions is Boosting, thanks to its good values of RPD in both calibration and validation terms ($1.6 \leq RPD_{cal, val} \leq 2.5$). Also in this case, RPD for Silt variable have unacceptable values ($RPD_{cal, val} \leq 1.2$).

Another good model in prediction of content of soils are RF, in fact in this case study, it has performed very good values of RPD both in calibration and validation sets for Clay, Sand and OC ($2 \leq RPD_{cal, val} \leq 2.3$), while fair and poor ones for Silt

($1.2 \leq \text{RPD}_{\text{cal, val}} \leq 1.5$). It is also very useful to highlight that RF has the best compromise, in RPD outcomes, to predict Silt among the various models performed.

Bagging is an alternative model that has carried out good RPD values between 1.5 and 2 for Clay, Sand, and OC, while between 1 and 1.5 for Silt. ANN instead, has performed lower values of RPD in Clay Sand and Silt, denoting itself as a poor model to predict soil texture but an excellent model in order to predict OC with values extremely good ($\text{RPD}_{\text{cal}} = 3.4$ and $\text{RPD}_{\text{val}} = 2.1$). The worst model performed is RT, which in all variables considered, carried out poor/slightly fair values of RPD in both calibration and validation datasets, but its usage is still not recommended.

5. Conclusions

This paper aims to evaluate the goodness of different multivariate and statistical ensemble methods in order to predict some soil properties.

Analogies and differences with our results appear in other papers, where authors applied similar techniques in different geographic areas, performing statistical calibration. For all properties we analyzed, PLSR is the technique that gave best results. This technique is the most complete and it is useful to predict many soil properties. Anyway, other alternative methods give good results. It would be worth if these techniques would be applied to deepen studies in soil properties predictions. Please refer to future studies in order to develop and broaden these issues, which are the subject of numerous papers.

Statement

All authors reviewed and revised the manuscript, approved the final version, and agreed to submit the revised manuscript for publication. The authors state that they have no disclosure to declare.

References

- AA.VV., (1995), *S-Plus, Guide to Statistics*, Seattle-WA: MathSoft.
- Akaike, H., (1973), 'Second International Symposium on Information Theory', in *Information theory and an extension of the maximum likelihood principle*, B. Petrov & F. Csaki, eds. Budapest: Akademiai Kiado, pp 267-281.

- Araújo, S.R., Wetterlind, J., Demattê, J.A.M., Stenberg, B., (2014), 'Improving the prediction performance of a large tropical vis-NIR spectroscopic soil library from Brazil by clustering into smaller subsets or use of data mining calibration techniques', *Eur. J. Soil Sci.* 65 (5), pp 718–729.
- Bellman, R., (1954) 'The theory of dynamic programming'. *Bull. Amer. Math. Soc.*, 60(6), pp. 503-515.
- Breiman, L., Friedman, J., Olshen, R. & Stone, C., (1984), *Classification and Regression Trees*, New York - London: Chapman & Hall.
- Brown, D.J., Shepherd, K.D., Walsh, M.G., Mays, M.D., Reinsch, T.G., (2006), 'Global soil characterization with VNIR diffuse reflectance spectroscopy', *Geoderma* 132 (3–4), pp273–290.
- Burger, J.; Geladi, P., (2007), 'Spectral pre-treatments of hyperspectral near infrared images: Analysis of diffuse reflectance scattering', *J. Near Infrared Spec* 15, pp 29–37.
- Clark, R.N.; Roush, T.L., (1984), 'Reflectance spectroscopy: Quantitative analysis techniques for remote sensing applications', *J. Geophys. Res.* 89, pp 6329–6340.
- Cozzolino, D. and Moròn A., (2003), 'The potential of near-infrared reflectance spectroscopy to analyse soil chemical and physical characteristics', *Journal of Agricultural Science*, 140, pp 65–71.
- Daubechies, I., (1992), 'Ten Lectures on Wavelets', *Society for Industrial and Applied Mathematics: Philadelphia, PA, USA*, p341.
- De Jong, S., (1992), 'The analysis of spectroscopical data to map soil types and soil crusts of Mediterranean eroded soils', *Soil Technol.* 5, pp 199–211.
- Di Gennaro, A., (2002), 'I sistemi di terra della Campania', *SELCA: Florence, Italy*, p 63.
- Ding, J.; Yang, A.; Wang, J.; Sagan, V.; Yu, D., (2018), 'Machine-learning-based quantitative estimation of soil organic carbon content by VIS/NIR spectroscopy', *Peer J.* 5714, pp 1–14.
- Drury, S.A.; (1993), *Image interpretation in geology*, Chapman & Hall: London, 1993.
- Dunn, B.W.; Beecher, H.G.; Batten, G.D.; Ciavarella, S., (2002), 'The potential of near-infrared reflectance spectroscopy for soil analysis - A case study from the Riverine Plain of South-Eastern Australia'. *Aust. J. Exp. Agric.*, 42, pp607–614.
- Efron, B.; Tibshirani, R.J., (1993), *An Introduction of the Bootstrap*, 1st ed., Chapman and Hall: New York, NY, USA, p 436.
- Ergon, R., (2009), 'Re-interpretation of NIPALS results solves PLSR inconsistency problem', *J. Chemom.* 23/1, pp 72-75.
- Eriksson, L.; Johansson, E.; Kettaneh-Wold, N.; Trygg, J.; Wilström, C.; Wold, S., (2006), *Multi-and Megavariate Data Analysis; Part I-Basic principles and applications*, Umetrics Academy: Umeå, Sweden, p 425.
- FAO, Food and Agriculture Organization of United Nations, (2015), 'Status of the World's soil resources', Rome, pp 607.
- Friedman, J. H., (1991), *Multivariate adaptive regression splines*, *Annals of Statistics* 19, pp 1-67.

- Fystro, G., (2002), 'The prediction of C and N content and their potential mineralisation in heterogeneous soil samples using Vis-NIR spectroscopy and comparative methods', *Plant Soil* 246, pp 139–149.
- Gene, G., Heath, M. & Wahba, G., (1979), 'Generalized Cross-Validation as a Method for Choosing a Good Ridge Parameter', *Technometrics* 21, 2, pp 215–230.
- Geoportale Regione Campania, (accessed on 10 February 2019), *Sistema Informativo Territoriale della Regione Campania*. Available online: <https://sit2.regione.campania.it/content/carta-utilizzazione-agricola-dei-suoli>.
- Gholizadeh, A., Borůvka, L., Saberioon, M., Vašát, R., (2016), 'A memory-based learning approach as compared to other data mining algorithms for the prediction of soil texture using diffuse reflectance spectra', *Remote Sens.* 8 (4), 341.
- Graupe, D. and Lynn, J.W. (1969) 'Some aspects regarding mechanistic modelling of recognition and memory', *Cybernetics* 3, pp 119-141.
- Hecht-Nielsen, R., (1987) 'Counter propagation networks', *Applied Optics* 26, pp4979-4984
- Helland, I., (1998), 'On the structure of partial least square regression', *Commu. Stat* 17, pp 311-388.
- Hopfield, J.J., (1982) 'Neural networks and physical systems with emergent collective computational abilities', *Proceedings of the National Academy of Sciences* 79, pp 2554-2558
- Irons, J.R.; Weismiller, R.A.; Petersen, G.W., (1989), '*Theory and Applications of optical remote sensing*'; G. Asrar, Ed.; Wiley: New York, pp 66-106.
- James, G., Witten, D., Hastie, T. & Tibishirani, R., (2013), 'An Introduction to Statistical Learning with Application in R', *G. Casella, S. Fienberg & I. Olkin, eds. Springer Texts in Statistics. New York: Springer Science+Business Media*, pp 1-426.
- Keesstra, S.D., Bouma, J., Wallinga, J., Tittonell, P., Smith, P., Cerdà, A., Montanarella, L., Quinton, J.N., Pachepsky, Y., Van der Putten, W.H., Bardgett, R.D., Moolenaar, S., Mol G., Jansen, B., Fresco, L.O., (2016) 'The significance of soils and soil science towards realization of the United Nations Sustainable: Development Goals', *SOIL An interactive open-access journal of the European Geosciences Union, Copernicus Publications*, Vol. 2 N°2, pp.111-128.
- Kuang, B.; Tekin, Y. and Mouazen, A.M., (2015), 'Comparison between artificial neural network and partial least squares for on-line visible and near infrared spectroscopy measurement of soil organic carbon, pH and clay content', *Soil & Tillage Research Science Direct* 146, pp 243-252.
- Lee, R.J. (1959) 'Generalization of learning in a machine'. *Proc. 14th ACM National Meeting*.
- Lee, K.S., Lee, D.H., Sudduth, K.A., Chung, S.O., Kitchen, N.R., Drummond, S.T., (2009), 'Wavelength identification and diffuse reflectance estimation for surface and profile soil properties', *Trans. ASABE* 52 (3), pp 683–695.
- Leone, A.P., (2000), 'Spettrometria e valutazione della riflettanza spettrale dei suoli nel dominio ottico 400–2500 nm', *Riv. Ital. Telerilevamento* 9, pp 3–28.
- Leone, A.P.; Viscarra-Rossel, A.R.; Amenta, P.; Buondonno, A., (2012), 'Prediction of soil properties with PLSR and vis-NIR spectroscopy: Application to Mediterranean soils from Southern Italy'. *Curr. Anal. Chem.*, 8, pp 283–299.

- Leone, A.P.; Leone, G.; Leone, N.; Galeone, C.; Grilli, E.; Orefice, N.; Ancona, V., (2019), 'Capability of diffuse Reflectance Spectroscopy to predict soil water retention and related soil properties in an irrigated lowland district of southern Italy', *Water Science and Technology* 11, 1712.
- Li, B.; Morris, J.; Martin, E.B., (2002), 'Model selection for partial least squares regression', *Chemom. Intell. Lab. Syst.* 1, pp 79–89.
- Lucadamo, A.; Amenta, P.; Leone, N. (2020), 'Soil texture prediction via reduced k-means principal component multinomial regression', *Socio-Economic Planning Sciences* (in press).
- Lucadamo, A.; Leone, A., (2015), 'Principal Component Multinomial Regression and spectrometry to predict soil texture', *Journal of chemometrics* 29 (9), pp 514–520.
- Martens, H., (1985), 'Multivariate Calibration', Dr. techn. Thesis, Technical University of Norway.
- Martens, H.; Naes, T., (1989), *Multivariate Calibration*, John Wiley & S.: Chichester, UK.
- McCarty, G.W.; Reeves, J.B., III; Reeves, V.B.; Follett, R.F.; Kimble, J.M., (2002), 'Mid-infrared and near-infrared diffuse reflectance spectroscopy for soil carbon measurement', *Soil Sci. Soc. Am. J.*, 66, pp 640–646.
- Milton; E.J. (1987), 'Principles of field spectroscopy', *Int. J. Remote Sens.*, 12, pp1807-1827.
- MIPAF, Ministero delle Politiche Agricole e Forestali, (2000), *Metodi di Analisi Chimica del Suolo*; Franco Angeli: Milan, Italy, p 536.
- Mouazen, A.M., Kuang, B., De Baerdemaeker, J., Ramon, H., (2010), 'Comparison among principal component, partial least squares and back propagation neural network analyses for accuracy of measurement of selected soil properties with visible and near infrared spectroscopy', *Geoderma* 158 (1), pp 23–31.
- Mouazen, A.M.; Maleki, M.R.; Cockx, L.; Van Meirvenne, M.; Van Holm, L.H.J.; Merckx, R.; De Baerdemaeker, J.; Ramon, H., (2009), 'Optimum three-point linkage set up for improving the quality of soil spectra and the accuracy of soil phosphorus measured using an on-line visible and near infrared sensor', *Soil and Tillage Research*, Volume 103, Issue 1, pp 144-152.
- Mouazen, A.M.; Maleki, M.R.; De Baerdemaeker, J.; Ramon, H., (2007), 'On-line measurement of some selected soil properties using a VIS-NIR sensor', *ScienceDirect Soil & Tillage Research* 93, pp 13-27.
- Mouazen, A.M.; Maleki, M.R.; De Ketelaere, D.; Ramon, H.; De Baerdemaeker, J.; (2008), 'On-the-go variable-rate phosphorus fertilisation based on a visible and near-infrared soil sensor', *Biosystems engineering Science direct* 99, pp 35-46.
- Nawar, S., Buddenbaum, H., Hill, J., Kozak, J., Mouazen, A.M., (2016), 'Estimating the soil clay content and organic matter by means of different calibration methods of vis-NIR diffuse reflectance spectroscopy', *Soil Tillage Res.* 155, pp 510–522.
- Pell, R., Ramos, L. & Manne, R., (2007), 'The model space in partial least squares regression', *J. Chemom.* 21, pp 165-172.
- Phatak, A. & De Jong, S., (1997), 'The geometry of partial least squares', *J. Chemo.* 11, pp 311-338.

- R Core Team, (2020), *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.
- Ripley, B., (1996), *Pattern Recognition and Neural Networks*. Cambridge: Univ Press.
- Rosenblatt, F., (1958) The perceptron, a probabilistic model for information storage and organization in the brain, *Psychol. Rev.* 65, pp 386-408.
- Rumelhart, D.E., Hinton, G.E. and Williams, R.J. (1986) 'Learning internal representations by error propagation', in *Parallel Distributed Processing. Explorations in the Microstructures of Cognition*, eds. Rumelhart, D.E. and McClelland, J.L. MIT press, Cambridge, MA.
- Shi, Z., Ji, W., Viscarra-Rossel, R.A., Chena, S., Zhou, Y., (2015), 'Prediction of soil organic matter using a spatially constrained local partial least squares regression and the Chinese vis-NIR spectral library', *European Journal of Soil Science*.
- Stenberg, B.; Viscarra-Rossel, R.A.; Mouazen, A.M.; Wetterlind, J. (2010), 'Visible and near infrared spectroscopy in soils science'. *Advances in Agronomy*, 107, pp 163-215.
- Stevens, A., Udelhoven T., Denis A., Tychon, B., Liroy R., Hoffmann L., Van Wesemael B., (2010), 'Measuring soil organic carbon in croplands at regional scale using airborne imaging spectroscopy', *Geoderma* 158, pp 32-45.
- Stevens, A.; Nocita, M.; Tóth, G.; Montanarella, L.; Van Wesemael, B.; (2013), 'Prediction of Soil Organic Carbon at the European Scale by Visible and Near InfraRed Reflectance Spectroscopy', *PLSoS One* 8(6), e66409.
- Vågen, T.G.; Shepherd, K.D.; Walsh, M.G., (2006), 'Sensing landscape level change in soil fertility following deforestation and conversion in the highlands of Madagascar using Vis-NIR spectroscopy', *Geoderma* 133, pp 281-294.
- Van der Meer, F., (1999), 'Can we map swelling clays with remote sensing?', *Int. J. Appl. Earth Obs. Geoinf.* 1, pp 27-35.
- Vasques, G.M.; Grunwald, S.; Sickman, J.O., (2008), 'Comparison of multivariate methods for inferential modelling of soil carbon using visible/near-infrared spectra'. *Geoderma*, 146, pp 14-25.
- Venables, W. & Ripley, B., (1997), '*Modern Applied Statistics with S-Plus*', Springer.
- Vibhute, A.D; Kale, K.V.; Mehrotra, S.C.; Dhumal, R.K.; Nagne, A.D., (2018), 'Determination of soil physicochemical in farming sites through visible, near-infrared diffuse reflectance spectroscopy and PLSR modeling', *Ecological Processes, Springer Open*, pp 7-26.
- Viscarra-Rossel, R.A., Walvoort, D.J.J., McBratney, A.B., Janik, L.J., Skjemstad, J.O.; (2006), 'Visible, near infrared, mid infrared or combined diffuse reflectance spectroscopy for simultaneous assessment of various soil properties', *Geoderma*, 131 (1), pp 59-75.
- Viscarra-Rossel, R.A., (2007), 'Robust modelling of soil diffuse reflectance spectra by bagging-partial least square regression', *J. Near Infrared Spec* 15, pp 39-47.
- Viscarra-Rossel, R.A., (2008), 'ParLeS: Software for chemometric analysis of spectroscopic data', *Chemom. Intell. Lab. Syst.* 90, pp72-83.
- Viscarra-Rossel, R.A., Lark, R.M., (2009), 'Improved analysis and modelling of soil diffuse reflectance spectra using wavelets', *Eur. J. Soil Sci.* 60 (3), pp 453-464.

- Viscarra-Rossel, R.A., Behrens, T., Ben-Dor, E., Brown, D.J., Demattê, J.A.M., Shepherd, K.D., Shi, Z., Stenberg, B., Stevens, A., Adamchuk, V., Aichi, H., Barthes, B.G., Baratholomeus, H.M., Bayer, A.D., Bernoux, M., Bottcher, K., Brodsky, L., Du, C.W., Chappell, A., Fouad, Y., Genot, V., Gomez, C., Grunwald, S., Gubler, A., Guerrero, C., Hedley, C.B., Knadel, M., Morras, H.J.M., Nocita, M., Ramirez-Lopez, L., Roudier, P., Rufasto Campos, E.M., Sanborn, P., Sellitto, V.M., Sudduth, K.A., Rawlins, B.G., Walter, C., Winowiecki, L.A., Hong, S.Y., Ji, W., (2016), 'A global spectral library to characterize the world's soil', *Earth Sci. Rev.* 155, pp 198–230.
- Vohland, M., Besold J., Hill J., Fründ H.C., (2018), 'Comparing different multivariate calibration methods for the determination of soil organic carbon pools with visible to near infrared spectroscopy', *Geoderma* 166, pp 198-205.
- Wang, S.Q., Li, W.D., Li, J., Liu, X.S., (2013), 'Prediction of soil texture using FT-NIR spectroscopy and PXRF spectrometry with data fusion', *Soil Sci.* 178 (11), pp 626–638.
- Were K., Tien Bui D., Dick Ø.B., Singh B.R., (2015), 'A comparative assessment of support vector regression, artificial neural networks, and random forests for predicting and mapping soil organic carbon stocks across an Afrotropical landscape', *Ecological Indicators* 52, pp 394-403.
- Wetterlind, J., Stenberg, B. and Jonsson, A., (2008), 'Near infrared reflectance spectroscopy compared with soil clay and organic matter content for estimating within-field variation in N uptake in cereals', *Plant and soil.* 302, 1-2, pp 317-327.
- Widrow, B. and Hoff, M.E., (1960) 'Adaptive switching circuits', *Proc. IRE WESCON conf.*, New York, pp. 96-104
- Widrow, B. and Winter, R. (1988) 'Neural nets for adaptive filtering and adaptive pattern recognition', *Computer*, 21, pp. 25-39.
- Williams, P.C., (1987), 'Variables affecting near-infrared reflectance spectroscopic analysis', *Near-Infrared Technology in the Agricultural and Food Industries*, American Association of Cereal Chemists Inc, pp 143-167.
- Wise, B.M.; Gallagher, N.B.; Bro, R.; Shaver, J.M., (2003), *PLS Toolbox Version 3.0 for Use with Matlab*, Eigenvector Research Inc., p 171.
- Wold, H., (1973), 'Non-linear iterative partial least squares (NIPALS) modelling. Some current developments'. In: Krishnaiah, P. R. (ed). *Multivariate Analysis III*, Academic Press, New York, pp. 383-407.
- Wold, S., Martens, H. & Wold, H., (1983), 'Lecture Notes in Mathematics'. In: A. Ruhe & B. Kagstrom, eds. *Proceedings of the Conference on Matrix Pencils*. Heidelberg, Germany: Springer-Verlag, pp 286-290.
- Wold, S., Johansson, E. & Cocchi, M., (1993) '3D QSAR in Drug Design, Theory, Methods, and Applications', *ESCOM Science Publishers: The Netherlands*, pp 523-550.
- Wold, S., Eriksson, L., Trygg, J. & Kettaneh, N., (2004), *Proceedings of COMPSTAT 2004*, Prague, Physica Verlag: Germany, pp 522-529.
- Xu S., Zhao Y., Wang M., Shi X., (2018), 'Comparison of multivariate methods for estimating selected soil properties from intact soil cores of paddy fields by Vis-NIR spectroscopy', *Geoderma* 310, pp 29-43.

Studio longitudinale sulle variazioni del (di)stress tra il pre e post-intervento chirurgico per carcinoma mammario in funzione delle strategie di coping adoperate

Maria De Caro¹, Ilaria Pepe^{1*}, Paolo Taurisano¹, Ernesto Toma²

¹ Dipartimento di Scienze mediche di base, neuroscienze e organi di senso,

² Dipartimento di Economia e Finanza

Riassunto: La letteratura scientifica ha proposto notevoli evidenze circa il ruolo dei fattori stressanti e delle strategie di coping nelle diverse fasi della malattia oncologica. Sulla base della teoria secondo cui il rischio di aumento/riduzione di stress psicologico sia altamente correlato con modalità di coping inadeguate/adequate di affrontare la patologia, il presente lavoro longitudinale si è posto l'obiettivo di dimostrare come le variazioni del distress tra il pre e post intervento chirurgico per cancro al seno fossero funzione delle strategie di coping adoperate. Il campione, reclutato presso il Centro di Senologia multidisciplinare del Policlinico di Bari, era costituito da 135 donne con diagnosi di neoplasia mammaria maligna. Oltre alla acquisizione tramite colloquio delle informazioni sociodemografiche e cliniche, sono stati somministrati il termometro del Distress per la rilevazione dello stress psicologico e la versione ridotta del Brief Cope per l'identificazione delle strategie di coping più utilizzate esclusivamente nel pre-operatorio. Le metodologie multivariate hanno evidenziato come le modalità di coping adattive fossero predittori statisticamente significativi della variabile dipendente, variazione del distress, perché ne hanno determinato l'abbassamento. Lo studio evidenzia l'importanza clinica delle risorse di coping nel rilevare le oscillazioni dello stress.

Keywords: Carcinoma mammario, distress, strategie di coping, intervento chirurgico.

* Autore corrispondente: i.pepe4@studenti.uniba.it

1. Introduzione

Come una tempesta, uno shock psichico che spinge le pazienti al di fuori del loro quieto mondo recondito, relazionale, sociale per giungere nel baratro della crisi e lacerazione interiore, l'evocazione del termine cancro vissuto come l'annuncio di un destino in declino, sommerge chi lo ascolta di un'ondata di sentimenti avversi, paure, angoscia, disperazione. Al termine cancro si sostituisce metaforicamente quello di una "fatal illness" impregnata di emotività e irrazionalità ed evocatrice di fantasie angosciose, divoranti, decostruenti. Il cancro al seno non è solo sviluppo proliferativo parassitico ma soprattutto invasione insidiosa delle componenti psichiche del paziente, indotte a sottoporre al vaglio della realtà interna ed esterna l'equivalenza cancro = lutto.

Il carcinoma mammario rappresenta la neoplasia invasiva più comune nelle donne, con circa 2 milioni di nuovi casi diagnosticati ogni anno nel mondo e un tasso di mortalità in crescita (627.000 morti solo nel 2018) È stimato come, almeno una donna su 8 nel corso della sua esistenza si ammali di tumore al seno e, soprattutto, tale probabilità aumenti in modo esponenziale dal momento in cui la donna rientra nella fase della menopausa (Who report on cancer, World Health Organization, 2020).

L'elevata incidenza di cancro al seno pone tale patologia come problematica da affrontare anche dal punto di vista sociale. Molti autori hanno indirizzato la loro attenzione sulle modalità attraverso le quali le diverse tipologie di interventi chirurgici (mastectomia radicale/chirurgia conservativa) inneschino risvolti psicologici in grado di incidere sulla qualità di vita delle pazienti (Carlson & Bultz, 2003; Mehnert & Koch, 2007).

Dal punto di vista psichico, infatti, la neoplasia mammaria è soprattutto sofferenza psicosociale, assume i connotati di evento traumatico in quanto minaccia l'integrità identitaria, l'autonomia personale, il senso di controllo, il tenore di vita. A causa delle reazioni emotive a breve e lungo termine, successive alla diagnosi, la membrana dell'ego corporeo, quale guscio protettivo, si strappa liberando pensieri intrusivi negativi. Il dolore incarnato nel corpo pone le donne dinanzi all'ombra della morte e alla percezione di finitudine, anche in coloro in grado di sopportare il peso di tale esperienza traumatica e di predisporre la strada verso una nuova trasformazione interiore accettando il nuovo corpo in tutte le sue funzioni. La natura cattiva e minacciosa della malattia sarà in ogni caso sperimentata come evocatrice di fantasmi di mutilazione, costituendo una vera e propria ferita alla femminilità.

Infatti, sebbene guariti dal cancro, le difficoltà relazionali, i deficit emotivi e quelli socio-familiari possono persistere (Conley et al., 2016; Bower, J.E. 2008).

Per questo motivo la storia anamnestica-emozionale delle pazienti costituisce l'elemento determinante nell'evoluzione del carcinoma mammario: ignorare questo legame significherebbe fallire nell'identificazione dei fattori di rischio che accrescono la probabilità di insorgenza di sintomatologia psichica. Infatti, nell'osteggiare la posizione di coloro che "portano in grembo" l'idea del cancro come un fenomeno di genetica molecolare, la letteratura scientifica ha proposto notevoli evidenze circa il ruolo della psiche e dei fattori stressanti nelle diverse fasi della malattia oncologica, dalla diagnosi agli interventi chirurgici, ai follow-up, dimostrando la relazione che sussiste tra la progressione della patologia e il sopraggiungere della sintomatologia stressante. I dati presenti in letteratura evidenziano che il 20- 45% delle donne affette da cancro al seno presenta sintomi di distress clinicamente significativi (Pugliese & Falcicchio, 2015). È possibile considerare il distress lungo un continuum: ad un estremo si collocano le ordinarie sensazioni di vulnerabilità, tristezza, paura, sperimentate di fronte all'esperienza di malattia, e all'altro si insedia una condizione di forte disagio che può manifestarsi con sintomi di ansia e di depressione, panico, isolamento sociale (Holland & Alici, 2010). Lo stato stressogeno presenta un andamento piuttosto variabile in oncologia, in letteratura sono stati riscontrati risultati contrastanti circa l'aumento/diminuzione di stress tra due fasi cruciali dell'intero iter-oncologico, quali il periodo precedente e successivo l'intervento chirurgico. Nonostante la maggior parte delle indagini abbia rilevato un'elevata incidenza al momento della scoperta e nella fase che precede l'asportazione, incidenza che decresce nel corso del tempo e a seguito dell'intervento, sebbene non ritorni mai ai normali livelli della popolazione generale, sul versante opposto, altri studi, hanno rivelato come a livello psicologico, il periodo successivo all'operazione, rappresenti un momento di stress acuto, in cui le pazienti si trovano a dover fronteggiare difficoltà fisiche, pratiche e psicologiche (Veronesi et al., 2002). Dunque, mentre le prime survey sono state in grado di dimostrare come i livelli di stress psicologico siano più elevati subito dopo la diagnosi e prima di sottoporsi all'intervento (Oh S, Miyamoto H et al., 2007) perché la preoccupazione principale è legata al timore della morte e al terrore dell'operazione (Piot-Ziegler et al., 2010), le seconde, invece, hanno evidenziato un accrescimento nella fase post-operatoria spiegabile in virtù del fatto che il distress psicologico sia direttamente correlato ad una disturbata immagine corporea: l'intervento chirurgico modifica una parte del corpo simbolo di tre tematiche fondamentali: la femminilità, la maternità e l'intimità.

L'adattamento alla malattia oncologica si configura anch'esso come un processo che si dispiega nel tempo e che va dalle precoci preoccupazioni vissute durante i primi sintomi, alla diagnosi, all'intervento chirurgico, al trattamento, alla remissione dalla malattia, alla preoccupazione per nuove recidive. Vos e colleghi (2004) hanno scoperto che nel periodo successivo l'intervento chirurgico, lo stile di coping rappresenta un parametro rilevante per instaurare un processo di adattamento psicosociale e per innescare risposte di fronteggiamento e conseguente alleviamento del carico stressante (Glinder et al., 2007).

Secondo Lazarus e Folkman (1984) il coping si riferisce al complesso delle strategie emotive, cognitive e comportamentali messe in atto per fronteggiare vissuti deleteri superiori alle risorse del soggetto, al fine di ridurre il loro impatto e garantirne una gestione più efficace. Nello specifico, nell'investigare il ruolo delle risposte allo stress a seguito di una diagnosi di cancro al seno è doveroso distinguere tra una forma di coping focalizzata sul problema (problem-focused) ed una centrata sull'emozione (emotion-focused) (Oswiecki D.M. & Compas B.E. 1999). Il ricorso all'una piuttosto che all'altra non discende tanto da stabili differenze di personalità, quanto dalla valutazione della situazione stressante esperita e dalla disponibilità di risorse individuali e relazionali.

La prima tipologia di coping racchiude attività tese a modificare, evitare o minimizzare l'urto dello stressor o, attraverso strategie cognitive, autoconvincersi di poter controllare lo stressor stesso. La seconda tipologia di coping viene generalmente considerata meno adattiva rispetto alle risposte centrate sui problemi in quanto fa riferimento a tentativi di dominare l'emozione attraverso la messa in atto di meccanismi che impediscano il conflitto diretto con l'evento stressante.

Poiché tale suddivisione proposta da Lazarus & Folkman non è stata considerata del tutto esaustiva, si è giunti ad una categorizzazione che comprende strategie cosiddette di "approach coping" versus "avoidant coping". La prima categoria attiene a quei comportamenti orientati al problema e dunque fa riferimento all'utilizzo di strategie sia cognitive che emotive nelle quali lo scopo è quello di ridurre e gestire le richieste interne o esterne derivanti dalla fonte stressante: ciò comporta trasformazioni dell'ambiente, del soggetto o della interazione tra i due. La seconda modalità concerne tutti quegli atteggiamenti di allontanamento dello stressor volti a ignorare ed evitare le conseguenze emotive dell'evento stressante: la negazione della minaccia e l'abbassamento della tensione mediante l'espressione emotiva sono un esempio (Aspinwall & Taylor, 1997; Endler & Parker, 1990; Suls & Flechter, 1985).

In particolare, è stato rilevato come, mentre l'attuazione di uno stile approach coping, sia esso cognitivo o comportamentale, correli con un minor livello di angoscia, abbassamento dei livelli di stress, migliore qualità di vita e maggiore possibilità di sopravvivenza (Gilbar et al., 2005; Montgomery et al., 2003; Stanton et al., 2002), le strategie di evitamento (avoidance coping), invece, fungono da predittori di elevato distress e peggiore conduzione dell'esistenza (Arraras, et al., 2002; Ben-Zur et al., 2001).

In tal senso, nell'ipotesi secondo cui le stimolazioni emotivamente disfunzionali rese note in due momenti diversi del percorso psiconcologico, pre e post-intervento chirurgico, siano strettamente dipendenti dalla capacità di mettere in atto modalità funzionali di adattamento alla situazione vissuta come traumatica, nel presente studio le strategie di coping si sono ritenute essere un fattore predittivo delle variazioni dello stress. Mentre tentativi rigidi di evitare/negare stati di attivazione fisiologica cronica possono produrre l'effetto di mantenere o incrementare maggiormente i livelli di stress psicologico, un atteggiamento attivo, combattivo, propositivo, orientato alla ricerca di soluzioni sarebbe, invece, in grado di ridurre la condizione di disagio catalizzata (Groarke et al., 2013; Hodges K & Winstanley S, 2012).

2. Obiettivi e ipotesi di ricerca

Il presente lavoro di ricerca ha adottato un disegno longitudinale al fine di esaminare se in un campione di donne affette da carcinoma mammario vi fossero differenze statisticamente significative del distress misurato in fase pre e post-chirurgica (Meijer A, Roseman M, Delisle VC, et al.2013), variazione che si è supposta essere in stretta correlazione con l'influenza esercitata dalle modalità più o meno funzionali di affrontare quel determinato momento storico della patologia (Arraras et al., 2002; Ben-Zur et al., 2001).

Sulla base della letteratura sopracitata, ex ante ci si aspettava che:

- H1. nel passaggio dal tempo t_0 al tempo t_1 , vi sarebbe stata una riduzione della quantità di distress;
- H2. le strategie disadattive (di evitamento) saranno positivamente associate all'aumento di distress psicologico, ergo, vi sarà un suo mantenimento/accrescimento;
- H3. le strategie di coping adattive (orientate sul problema) saranno negativamente correlate alla diminuzione del distress, innescandone l'abbassamento.

3. Materiali e metodi

3.1 Partecipanti

Sono state reclutate per lo studio 135 donne, sostenute e curate mediante un intervento psicologico effettuato durante il lasso di tempo che intercorre tra la prima visita precedente alla rimozione di neoplasia maligna alla mammella mediante terapia chirurgica e la visita successiva all'asportazione. Le pazienti sono state assistite da un gruppo di professionisti altamente specializzato nella diagnosi e cura del carcinoma alla mammella afferente al Centro di Senologia multidisciplinare del Policlinico di Bari. Sono state incluse donne con un'età superiore o equivalente a 18 anni affette da carcinoma mammario e operate presso il Policlinico di Bari Ospedale "Giovanni XXIII" (Padiglione Balestrazzi e Padiglione Asclepios) dal team della Breast Care Unit. Tra queste, l'intero campione ha accettato di sottoporsi al protocollo diagnostico e terapeutico proposto dall'equipe multidisciplinare nonostante una parte avesse effettuato la diagnosi esternamente. Non sono state campionate pazienti che, in anamnesi o al momento della valutazione, presentassero patologie psichiatriche categorizzate nel sistema nosografico più utilizzato a livello internazionale (DSM-5) o soddisfacessero criteri di comorbidità con altre condizioni mediche generali, ad esclusione della neoplasia maligna. Anche l'abuso o dipendenza da sostanze, corrente o rilevato nella storia anamnestica, ha rappresentato un parametro di esclusione.

3.2 Strumenti

A seguito di un colloquio psicologico predisposto per indagare le seguenti aree: dati socio-anagrafici, prima diagnosi/recidiva e familiarità della malattia ovvero raggruppamento di più casi di cancro al seno nel nucleo familiare senza un'evidente trasmissione della malattia da una generazione alla successiva, sia durante la fase pre-chirurgica sia nel lasso di tempo successivo alla rimozione della neoplasia maligna della mammella, la valutazione psicometrica del distress psicologico è stata eseguita mediante uno strumento chiamato *Termometro del Distress*. Tale test di screening, costituito da un singolo item, valuta lo stress soggettivo percepito attraverso una scala analogica visiva (a forma di termometro) rispetto alla quale le pazienti devono attribuire un valore che va da 0 (assenza di disagio) a 10 (disagio elevato) e che meglio descrive il livello di stress vissuto nell'ultima settimana. Il cut-off ideale per individuare la presenza di distress psicologico moderato nei pazienti oncologici è rappresentato da un punteggio compreso tra 4 e 7, mentre un punteggio superiore a 7 suggerisce un livello di distress francamente patologico (Gil et al., 2005).

Le strategie di coping, invece, sono state misurate nel pre-intervento utilizzando la versione italiana del *Brief Cope- dispositional version* (Carver, 1997; Conti, 1999), versione ridotta del Cope Inventory (Carver et al., 1989). Esso è concepito come uno strumento di misura in grado di identificare più sottili differenze individuali di coping, dimostrandosi capace di bilanciare la propensione generale del soggetto (“in che modo reagirebbe se”) con la risposta corrente alla situazione stressante. Sostanzialmente si tratta di un questionario self-report costituito da 28 item e progettato per misurare modalità efficaci/inefficaci di far fronte ad un evento di vita stressante. Ai partecipanti viene chiesto di rispondere valutando ciascun item su una scala a 4 punti, da 1 (mai) a 4 (sempre). La scala viene spesso utilizzata in ambito sanitario al fine di valutare in che modo i pazienti stiano affrontando la fase diagnostica di una particolare patologia invasiva. Essa comprende 14 sottoscale che corrispondono ad altrettante diverse strategie di coping, di cui 6 attinenti a stili di avoidant coping e 6 a strategie di approach coping.

3.3 Metodologie statistiche

Tutte le analisi statistiche sono state effettuate mediante il software SPSS versione 20, utilizzando, per tutti i test applicati, un livello di significatività $\alpha=0,05$. Sono state eseguite analisi descrittive sulle caratteristiche del campione e sugli indici TD e Brief Cope espresse con media e deviazione standard.

Al fine di verificare se vi fossero variazioni statisticamente significative nei livelli di distress rilevati in due momenti differenti: periodo pre-operatorio (t_0) e fase post-chirurgica (t_1), è stata condotta una verifica di ipotesi sull'uguaglianza tra due medie nel caso di dati appaiati (i soggetti non differiscono nei due tempi di rilevazione). È stato utilizzato il test non parametrico di Wilcoxon per due campioni dipendenti dopo aver accertato l'assenza di normalità della variabile “DeltaDistress” (differenza tra i punteggi tra il tempo 0 e il tempo 1), requisito indispensabile per la corretta applicazione del test T su campioni dipendenti.

Con l'obiettivo di appurare la misura in cui la valenza adattiva o meno dello stile di coping dei pazienti, adottato nel periodo pre-operatorio, possa predire le variazioni del distress tra la fase precedente all'intervento chirurgico e quella successiva, ergo, confermare il ruolo fondamentale assunto dalle strategie di coping nello spiegare la diminuzione/incremento della variabile DeltaDistress, è stata eseguita un'analisi di regressione lineare multipla. Prima di giungere all'applicazione di tale tecnica multivariata, con la finalità di riassumere e semplificare le relazioni dell'insieme di modalità di coping adoperate, è stata utilizzata una analisi fattoriale. Successivamente all'estrapolazione dei fattori concernenti le differenti strategie di coping, si è giunti

ad applicare un metodo di regressione lineare multipla tra la variabile dipendente (VD) DeltaDistress e le variabili indipendenti (VI) rappresentate dalle componenti fattoriali di cui sopra. L'ipotesi nulla prevede che non sussista alcun tipo di associazione tra VD e VI, mentre l'ipotesi alternativa sostiene che almeno uno dei predittori, in particolar modo quello legato all'utilizzo di un cospicuo numero di modalità di coping funzionali, contribuirà a spiegare significativamente le variazioni della VD.

4. Risultati

4.1 Caratteristiche socio-demografiche del campione

Il campione, costituito da 135 donne con diagnosi di neoplasia maligna della mammella e reclutato nello studio ha un'età media di 59,87 anni (Std. Dev.=12,44, range 28-82) (Tab.1). Sono state individuate quattro classi caratterizzate da distribuzione numerica omogenea: poco meno di 3/4 del campione ha meno di 70 anni e tutte le singole età da 41 a 82 (valore Max) sono presenti con poche unità statistiche (Tab.1bis). Al di sotto dei 41 anni rientrano solo tre soggetti (Fig.1).

Una certa equidistribuzione campionaria è riscontrabile anche in relazione alla variabile scolarità: tra i rispondenti, nonostante la maggior parte abbia conseguito il diploma di scuola media superiore (29,1%), un non indifferente 22,4% è fermo alla licenza elementare, percentuale analoga di coloro che, invece, sul versante opposto, hanno ottenuto una laurea di I o II livello (Tab.2).

Per quanto concerne l'attività lavorativa, è stato riscontrato come, ad eccezione di tre soggetti che non hanno offerto alcun tipo di risposta, più della metà lavori (61,4%) (Tab.3) e tra il totale delle occupate (n=81), è la mansione dell'impiegata ad essere svolta maggiormente (29,6%), seguita dalla professione dell'insegnante (21%) (Tab.3bis).

Esaminando il contesto familiare, è possibile constatare che, la più cospicua parte del campione, escluse le non rispondenti, è coniugata (75,2%) (Tab.4) e ha mediamente 2 figli (Tab.5bis).

In riferimento alle variabili cliniche, ed in particolar modo a due momenti diversi del percorso psiconcologico: la prima diagnosi e la recidiva, le partecipanti selezionate per lo studio in prima diagnosi sono pari all'80,7%, con solo il 19,3% di pazienti con recidiva (Tab.6). Infine, considerevole la percentuale dei soggetti (61,2%) (Tab.7) che è possibile definire "familiare" per la presenza di alcuni casi della stessa malattia in un singolo ambito familiare ma non in rapporto diretto.

Tabella 1. Statistiche descrittive in relazione alla variabile età

	N	Minimo	Massimo	Media	Std. Deviation
Età (in anni compiuti)	135	28	82	59,87	12,44

Tabella 1 bis. Distribuzione, in v.a. e %, delle pazienti, per classi d'età (in anni compiuti)

Classi di età (in anni compiuti)	Frequenze assolute	%	Frequenze cumulate
Meno di 50	30	22,2	22,2
50-59	36	26,7	48,9
60-69	33	24,4	73,3
70 e oltre	36	26,7	100,0
Totale	135	100,0	

Figura 1. Distribuzione in v.a. delle pazienti per età (in anni compiuti)

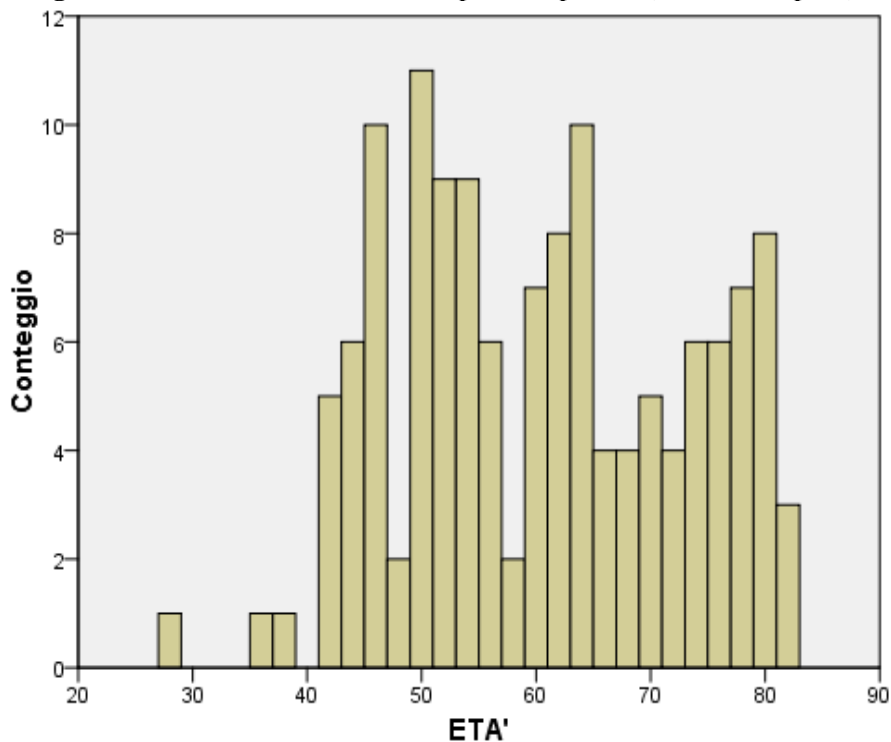


Tabella 2. Distribuzione, in v.a. e %, delle pazienti, per scolarità

Titolo massimo di studio conseguito	Frequenze assolute	%	% valide	Frequenze cumulate
Licenza elementare	30	22,2	22,4	22,4
Licenza media inf.	35	25,9	26,1	48,5
Licenza media sup.	39	28,9	29,1	77,6
Laurea (I o II livello)	30	22,2	22,4	100,0
Totale	134	99,3	100,0	
Valori mancanti	1	,7		
Totale	135	100,0		

Tabella 3. Distribuzione, in v.a. e %, delle pazienti, per condizione lavorativa e non

Condizione lavorativa e non	Frequenze assolute	%	% valide
Casalinga	51	37,8	38,6
Lavoratrice	81	60,0	61,4
Totale	132	97,8	100,0
Valori mancanti	3	2,2	
Totale	135	100,0	

Tabella 3bis. Distribuzione, in v.a. e %, delle pazienti, per professione svolta

Professione svolta	Frequenze assolute	%
Impiegata	24	29,6
Insegnante	17	21,0
OSS	5	6,2
Agricoltore	4	4,9
Imprenditrice	4	4,9
Parrucchiera-estetista	4	4,9
Sarta	4	4,9
Artigiana	2	2,5
Medico	2	2,5
Operaia	2	2,5
Architetto	1	1,2
Biologa	1	1,2
Commerciante	1	1,2
Consulente	1	1,2
Fisioterapia	1	1,2
Ingegnere	1	1,2
Musicista	1	1,2
Operatore ecologico	1	1,2
Ostetrica	1	1,2
Polizia penitenziaria	1	1,2
Psicoterapeuta	1	1,2
Tecnico di laboratorio	1	1,2
Traduttrice	1	1,2
Totale pazienti occupate	135	100,0

Tabella 4. Distribuzione, in v.a. e %, delle pazienti, per stato civile

Stato civile	Frequenze assolute	%	% valide
Coniugata	100	74,1	75,2
Nubile	10	7,4	7,5
Separata-Divorziata	8	5,9	6,0
Vedova	15	11,1	11,3
Totale	133	98,5	100,0
Valori mancanti	2	1,5	
Totale	135	100,0	

Tabella 5. Distribuzione, in v.a. e %, delle pazienti, per numero di figli

Numero di figli	Frequenze assolute	%	% valide	Frequenze cumulate
0	17	12,6	12,7	12,7
1	21	15,6	15,7	28,4
2	64	47,4	47,8	76,1
3	27	20,0	20,1	96,3
4	5	3,7	3,7	100,0
Totale	134	99,3	100,0	
Valori mancanti	1	,7		
Totale	135	100,0		

Tabella 5bis. Statistiche descrittive in relazione alla variabile numero di figli

	N	Minimo	Massimo	Media	Std. Deviation
Numero di figli	134	0	4	1,87	1,00

Valori mancanti: 1

Tabella 6. Distribuzione, in v.a. e %, delle pazienti, in prima diagnosi o recidiva

Prima diagnosi o recidiva	Frequenze assolute	%
Prima diagnosi	109	80,7
Recidiva	26	19,3
Totale pazienti	135	100,0

Tabella 7. Distribuzione, in v.a. e %, delle pazienti, per familiarità della patologia

Familiarità della patologia	Frequenze assolute	%	% valide
No	52	38,5	38,8
Sì	82	60,7	61,2
Totale	134	99,3	100,0
Valori mancanti	1	,7	
Totale	135	100,0	

4.2 Distress psicologico e coping: analisi descrittive delle misure

Le tabelle 8, 9 e 11 riportano la media dei punteggi della scala del coping (Brief COPE) somministrata nel pre-intervento, unico momento in cui il questionario viene raccolto, e dei livelli di distress valutati attraverso il Termometro del Distress (TD) in entrambe le misurazioni temporali ovvero prima di sottoporsi all'intervento chirurgico (t_0) e dopo averlo effettuato (t_1). In particolare, le analisi descrittive preliminari eseguite hanno mostrato come le pazienti valutate nel pre-operatorio presentassero un livello di distress mediamente moderato e tendente all'estremità più prossima alla definizione di un grado di disagio psicologico clinicamente rilevante ($M=6.61$, $DS=2,60$) (Tab.8), osservazione avvalorata dai valori pari a 7 assunti sia dalla mediana che dalla moda (Tab.10). In relazione alla variabile post-intervento, il punteggio medio, quello mediano e quello modale tendono ad assumere valori decisamente inferiori ($M=5.40$, $DS=2,54$; $Mediana=5$, $Moda=5$) (Tabb.9 e 10). Più nello specifico, per quanto riguarda la prevalenza di stress psicologico osservata al tempo 0, ben 63 pazienti su 135 (46,7%) hanno riportato un punteggio al TD compreso tra 4 e 7, indicativo di una quota di distress moderata. 54 donne (40,0%), invece, si assestavano attorno ad un livello sensibilmente superiore del cut-off pari a 8, sintomatico della presenza di stress clinicamente elevato (Tab.8bis). A conferma dell'abbassamento del punteggio medio ottenuto al TD nel frangente successivo all'asportazione del tessuto mammario, è possibile osservare che, a fronte della riduzione della numerosità campionaria nella manifestazione di una quota di distress severo e dell'aumento dei casi appartenenti alla categoria "assente"/"lieve", una considerevole parte delle pazienti, il 70,4%, oscilla tra la presenza di distress tenue e moderato e per questo collocabile tra i casi borderline, a rischio di evoluzione clinicamente rilevante (Tab.9bis).

I punteggi medi e le deviazioni standard dei 14 domini della scala del Brief-COPE sono riportati nella tabella 11. Le valutazioni medie riportate mostrano come il campione abbia utilizzato prevalentemente strategie di approach coping, in particolare "accettazione" ($M=5.59$, $DS=2,00$), "affrontare positivamente" ($M=5.31$, $DS=2,03$), "uso del supporto emotivo" ($M=5.02$, $DS=1,85$), "pianificazione" ($M=4.89$, $DS=2,18$), "uso del supporto strumentale" ($M=4.69$, $DS=1,85$), rispetto a quelle classificabili come avoidant coping ("uso di sostanze" con $M=2.11$ e $DS=,54$; "disimpegno comportamentale" con $M=2.61$ e $DS=1,39$; "autoaccusa" con $M=3.00$ e $DS=1,53$). Fa eccezione tra queste ultime il dominio "distogliere l'attenzione" che, con un punteggio medio di 5.45 ($DS=2,05$) si colloca al secondo posto nella gerarchia delle modalità di coping più frequentemente impiegate.

Tabella 8. Statistiche descrittive in relazione alla variabile *Distress pre-operatorio*

	N	Minimo	Massimo	Media	Std. Deviation
Distress pre-operatorio	135	0,00	10,00	6,61	2,60

Tabella 8bis. Distribuzione, in v.a. e %, delle pazienti, per livelli di Distress pre-operatorio

Livelli di Distress pre-operatorio	Frequenze assolute	%	Frequenze cumulate
ASSENTE	5	3,7	3,7
LIEVE	13	9,6	13,3
MODERATO	63	46,7	60,0
SEVERO	54	40,0	100,0
Totale	135	100,0	

Tabella 9. Statistiche descrittive in relazione alla variabile *Distress post-operatorio*

	N	Minimo	Massimo	Media	Std. Deviation
Distress post-operatorio	135	0,00	10,00	5,40	2,54

Tabella 9bis. Distribuzione, in v.a. e %, delle pazienti, per livelli di Distress post-operatorio

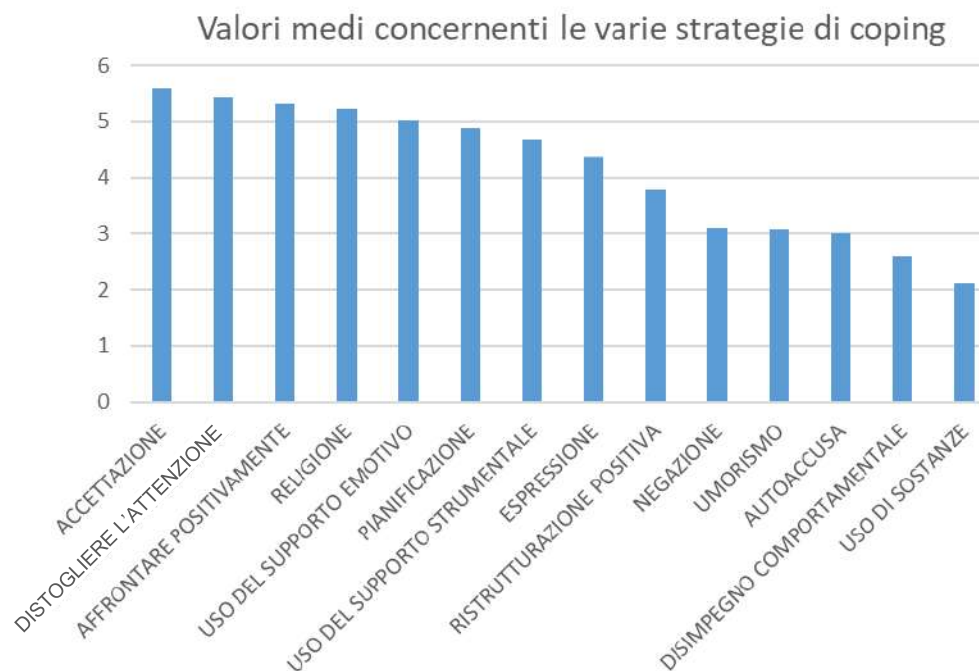
Livelli di Distress post-operatorio	Frequenze assolute	%	Frequenze cumulate
ASSENTE	9	6,7	6,7
LIEVE	14	10,4	17,0
MODERATO	81	60,0	77,0
SEVERO	31	23,0	100,0
Totale	135	100,0	

Tabella 10. Statistiche a confronto tra Distress pre e Distress post-operatorio

Statistiche	Distress pre-operatorio	Distress post-operatorio
Media	6,61	5,40
Mediana	7	5
Moda	7	5
Primo Quartile	5	4
Terzo quartile	8	7

Tabella 11. Statistiche relative alle 14 strategie di coping

Strategie di coping	Media	Std. Deviation
Espressione	4,37	1,85
Uso del supporto strumentale	4,69	1,85
Uso del supporto emotivo	5,02	1,85
Distogliere l'attenzione	5,45	2,05
Negazione	3,10	1,71
Umorismo	3,07	1,71
Disimpegno comportamentale	2,61	1,39
Uso di sostanze	2,11	,54
Ristrutturazione positiva	3,80	2,05
Accettazione	5,59	2,00
Affrontare positivamente	5,31	2,03
Pianificazione	4,89	2,18
Religione	5,21	2,10
Autoaccusa	3,00	1,53

Figura 2. Valori medi delle 14 strategie di coping

4.3 Verifica di ipotesi sull'uguaglianza del distress medio tra i due tempi di rilevazione

Non avendo potuto applicare il test T di Student per campioni dipendenti in quanto la variabile "DeltaDistress" (differenza tra i punteggi) di cui sono riportate le principali statistiche descrittive nella tabella 12, non si distribuisce in modo gaussiano (come confermato dal test di normalità di Shapiro- Wilk e da quello di Kolmogorov-Smirnov con la correzione di Lilliefors, entrambi con $p\text{-value} < \alpha = 0,05$, Tab. 13), attraverso il test di Wilcoxon per dati appaiati, si evince come, fissando un livello $\alpha = 0,05$, poiché $T_w = -5,08 < -z_{\alpha/2} = -1,96$ (statistica test standardizzata utilizzata essendo $n > 30$; Tab.14) l'ipotesi di base di uguaglianza delle due distribuzioni a confronto deve essere rifiutata, pertanto il livello di distress manifestato tra la fase pre-operatoria e quella post-operatoria è statisticamente differente, con una significativa diminuzione dello stesso nel passaggio al tempo 1 come già era stato constatato da un punto di vista meramente descrittivo nel precedente paragrafo (Tab.10) e confermato graficamente dalla Fig. 3. Quest'ultima mostra un decremento di tutte le statistiche (media, mediana, moda, primo e terzo quartile) nel passaggio dal periodo pre a quello post operatorio.

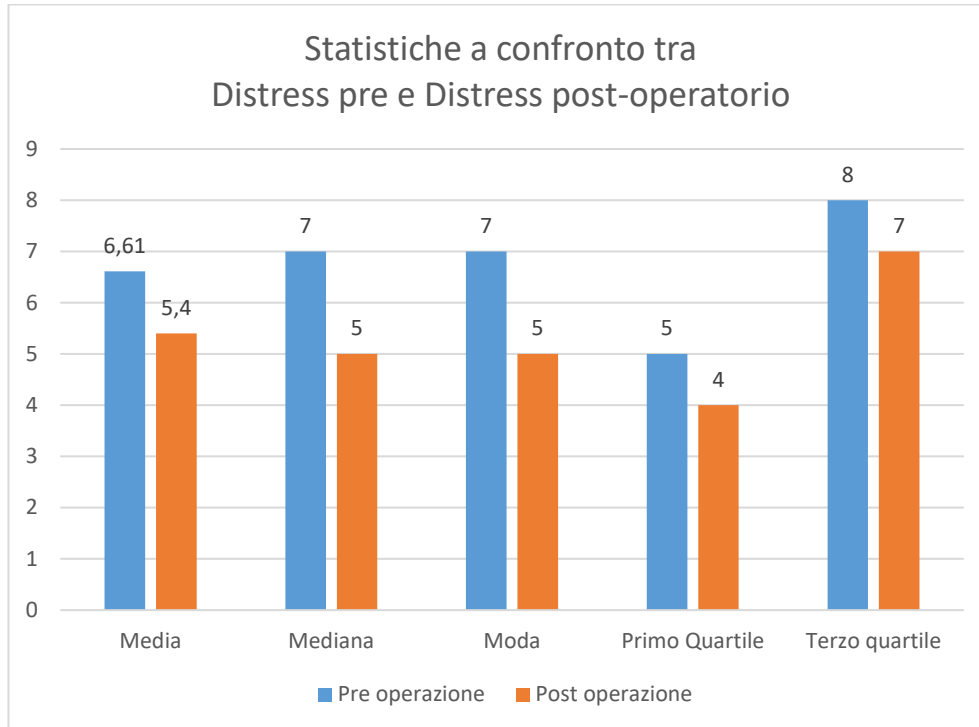
Tabella 12. Statistiche relative alla variabile DeltaDistress (Differenza tra Distress post-operatorio e Distress pre-operatorio)

Statistiche	Media	Std. Deviation
Media	-1,21	0,23
Intervallo di confidenza per la media al 95%		
<i>Limite inferiore</i>	-1,66	
<i>Limite superiore</i>	-0,75	
Mediana	0,00	
Varianza	7,09	
Deviazione std.	2,66	
Minimo	-10	
Massimo	5	
Intervallo	15	
Distanza interquartilica	2	
Asimmetria	-1,08	0,21
Curtosi	2,73	0,41

Tabella 13. Test di normalità relativi alla variabile DeltaDistress

Kolmogorov-Smirnov ^a			Shapiro-Wilk		
Statistica	Df	Sig.	Statistica	Df	Sig.
,21	135	0,00	,86	135	0,00

a. Correzione di significatività di Lilliefors

Figura 3. Statistiche a confronto tra Distress pre e Distress post-operatorio**Tabella 14.** Test di Wilcoxon sulla uguaglianza tra il Distress post-operatorio ed il Distress pre-operatorio nel caso di dati appaiati

Numero di casi totali	135
Statistica test	568,50
Errore standard	207,02
Statistica test standardizzata	-5,08
Sign.	0,00

4.4 Le strategie di coping predicono l'andamento del Distress? Analisi multivariate

L'Analisi Fattoriale ha consentito di partire dalle 14 dimensioni del Brief Cope e ottenerne 6 che rendono conto delle similarità che accomunano le variabili di partenza. I valori assunti dal test KMO e dal test di sfericità di Bartlett hanno permesso di effettuare l'analisi fattoriale perché soddisfatti i corrispettivi criteri richiesti quali campione sufficientemente ampio, test statisticamente significativo e un valore del KMO vicino a .70 (Tab.15).

Tabella 15. Test di adeguatezza di Kaiser-Meyer-Olkin e test di sfericità di Bartlett

Test di adeguatezza di Kaiser-Meyer-Olkin		0,66
Test di sfericità di Bartlett	χ^2 approssimato	497,83
	Gdl	91
	Sig.	0,00

Dopodiché, l'analisi fattoriale eseguita con il metodo delle componenti principali (ACP), come evinto dalle tabelle 16 e 16bis, ha consentito di individuare 6 dimensioni latenti che spiegano all'incirca il 72% della variabilità iniziale. Nello specifico, il primo fattore che rende conto del 24% della variabilità totale è correlato alla cospicua messa in atto di strategie di approach coping quali uso del supporto strumentale, uso del supporto emotivo, ristrutturazione positiva, accettazione, affrontare positivamente e pianificazione; la seconda, invece, che spiega circa il 13% della variabilità complessiva, è legata primariamente a due strategie disadattive, la negazione e il disimpegno comportamentale; la terza, che spiega poco più dell'11%, è costituita principalmente da modalità afferenti a categorie diverse tra loro: l'umorismo di stampo "neutrale", la ristrutturazione positiva di valenza adattiva e l'autoaccusa, modalità disfunzionale. Le restanti tre componenti individuate dall'analisi fattoriale spiegano tra il 7 e l'8% della variabilità, in particolar modo mentre la quarta è correlata negativamente con la strategia "distogliere l'attenzione", la quinta e la sesta sono associate all'uso di sostanze e alla religione.

Tabella 16. Autovalori e quota di variabilità spiegata dalle varie componenti estratte con l'Analisi fattoriale

Componenti	Autovalori iniziali			Pesi delle componenti estratte		
	Totale	% varianza	% cumulata	Totale	% varianza	% cumulata
1	3,38	24,16	24,16	3,38	24,16	24,16
2	1,81	12,94	37,11	1,81	12,94	37,11
3	1,59	11,38	48,49	1,59	11,38	48,49
4	1,18	8,42	56,91	1,18	8,42	56,91
5	1,08	7,69	64,60	1,08	7,69	64,60
6	1,01	7,22	71,82	1,01	7,22	71,82
7	0,84	6,00	77,82			
8	0,69	4,94	82,76			
9	0,63	4,50	87,26			
10-	0,51	3,62	90,88			
11	0,41	2,91	93,79			
12	0,39	2,81	96,61			
13	0,27	1,90	98,50			
14	0,21	1,50	100,00			

Tabella 16bis. Matrice di correlazione delle componenti principali estratte attraverso l'analisi fattoriale e le 14 strategie di coping

Strategie di coping	Componenti					
	1	2	3	4	5	6
Espressione	0,41	0,36	-0,36	0,19	-0,37	-0,09
Uso del supporto strumentale	0,73	0,07	-0,36	0,34	0,01	-0,01
Uso del supporto emotivo	0,70	0,03	-0,40	0,36	0,06	-0,17
Distogliere l'attenzione	0,50	0,19	-0,06	-0,65	-0,07	-0,10
Negazione	0,20	0,60	0,01	-0,22	0,33	0,18
Umorismo	0,40	0,30	0,51	0,21	-0,42	-0,07
Disimpegno comportamentale	0,29	0,73	-0,03	-0,12	0,03	0,16
Uso di sostanze	0,12	0,14	0,40	0,37	0,55	-0,39
Ristrutturazione positiva	0,51	-0,26	0,55	0,17	0,14	0,18
Accettazione	0,49	-0,42	0,30	0,00	-0,36	0,23
Affrontare positivamente	0,67	-0,29	0,01	-0,35	0,32	-0,22
Pianificazione	0,74	-0,39	-0,03	-0,26	-0,04	-0,16
Religione	0,25	-0,23	-0,33	0,13	0,30	0,71
Autoaccusa	0,30	0,34	0,50	-0,03	0,01	0,29

La matrice dei pesi fattoriali (Tab. 16bis) contiene valori piuttosto elevati che ben mettono in luce le relazioni tra le variabili originarie e i fattori comuni; nonostante ciò, attraverso le rotazioni dei fattori effettuate secondo i principali metodi in letteratura, si è provato a migliorare l'interpretazione degli stessi, non ottenendo, però, i risultati aspettati. Si presume, pertanto, che i fattori siano già ben distribuiti sugli assi indentificati dalle componenti principali.

A seguito dell'individuazione delle sei componenti principali, dall'effettuazione di una prima regressione, la tabella ANOVA (Tab.17) ha confermato la significatività dell'analisi della quota di varianza spiegata dalla relazione regressiva (Sig.=,02), permettendo di affermare come le variabili indipendenti del modello siano complessivamente predittori significativi della variabile dipendente. Tuttavia, dal momento che nella regressione lineare multipla i predittori sono più di uno, si rende necessario verificare quale sia il più importante in termini di potere euristico. Dunque, oltre alla bontà generale del modello, per conoscere la significatività di ciascun predittore, sono stati analizzati i coefficienti di regressione riportati nella tabella 17bis. La prima colonna della tabella contiene le sei variabili inserite nel modello, mentre nelle colonne successive è possibile leggere, per ciascun predittore, il valore del coefficiente di regressione (B) con il relativo errore standard, i coefficienti di regressione calcolati sulle variabili standardizzate (Beta) e il valore di t con relativa probabilità associata (Sig.). Poiché risultano essere statisticamente significative (p-value=,05; p-

value \leq 0,01, $\alpha=0.05$) esclusivamente la prima e la terza variabile predittiva, attraverso la tecnica della Backward deletion sono stati eliminati quei predittori che presentavano il coefficiente di regressione statisticamente meno significativo, giungendo ad un modello finale che ha confermato quanto visionato nella tabella 16bis.

Tabella 17. Anova relativa al modello regressivo tra Variabile dipendente: DeltaDistress (Differenza tra Distress post-operatorio e Distress pre-operatorio) e le sei variabili estratte dall'analisi fattoriale

Fonte di variabilità del modello	Somma quadrati	Gdl	Varianza	F	Sig.
Regressione	100,78	6	16,80	2,53	0,02 ^b
Residua	849,41	128	6,64		
Totale	950,19	134			

- a. Variabile dipendente: DeltaDistress (Differenza tra Distress post-operatorio e Distress pre-operatorio).
 b. Variabili indipendenti: (Constant), REGR factor score 6, REGR factor score 5, REGR factor score 4, REGR factor score 3, REGR factor score 2, REGR factor score 1 .

Tabella 17bis. Coefficienti del modello regressivo multiplo tra Variabile dipendente: DeltaDistress (Differenza tra Distress post-operatorio e Distress pre-operatorio) e le sei variabili estratte dall'analisi fattoriale

Variabili indipendenti	Coefficienti NON standardizzati		Coefficienti standardizzati	T	Sig.
	B	Std. Error	BETA		
Costante	-1,21	,22		-5,45	0,00
Componente fattoriale 1	-,45	,22	-,17	-2,02	0,05
Componente fattoriale 2	-,03	,22	-,01	-,12	0,90
Componente fattoriale 3	,59	,22	,22	2,65	0,01
Componente fattoriale 4	,41	,22	,15	1,86	0,06
Componente fattoriale 5	-,15	,22	-,05	-,66	0,51
Componente fattoriale 6	,10	,22	,04	,45	0,66

- a. Variabile dipendente: DeltaDistress (Differenza tra Distress post-operatorio e Distress pre-operatorio)

Rispetto all'ultimo modello a cui si è giunti, la tabella 18 dell'ANOVA dimostra come la variabilità dovuta alla regressione e legata alla relazione di dipendenza tra le due X (prima componente strutturata sulla base delle strategie di approach coping e terzo fattore concernente strategie tra loro dissimili) e la Y (DeltaDistress) sia statisticamente significativa (sig.<,01).

L'ANOVA di regressione, confrontando la componente di variabilità scaturita dalla regressione con la variabilità residua, ha permesso di giungere ad affermare che le variabili indipendenti siano predittori estremamente significativi della variabile

dipendente. Se questa assunzione corroborata permette di spiegare una generica predizione, più nello specifico la tabella 18bis dimostra come, mentre il primo fattore costituito da un consistente numero di modalità di approach coping sia inversamente proporzionale alla diminuzione del distress, la seconda componente, positivamente associata a tipologie di strategie tra loro differenti rispetto alla valenza attribuitagli, sia correlata positivamente con la variabile DeltaDistress.

Ai fini puramente descrittivi, va comunque sottolineato come la componente correlata con due particolari modalità di coping disadattive, nonostante la sua importanza in termini di variabilità nell'analisi fattoriale, non è risultata essere statisticamente associata alla variabile dipendente DeltaDistress.

Tabella 18. Anova relativa al modello regressivo tra Variabile dipendente: DeltaDistress (Differenza tra Distress post-operatorio e Distress pre-operatorio) e le due variabili statisticamente significative estratte dall'analisi fattoriale

Fonte di variabilità del modello	Somma quadrati	gdl	Varianza	F	Sig.
Regressione	73,55	2	36,78	5,54	0,00 ^b
Residua	876,63	132	6,64		
Totale	950,19	134			

- a. Variabile dipendente: DeltaDistress (Differenza tra Distress post-operatorio e Distress pre-operatorio)
 b. Variabili indipendenti: (Constant), REGR factor score 3, REGR factor score 1

Tabella 18bis. Coefficienti del modello regressivo multiplo tra Variabile dipendente: DeltaDistress (Differenza tra Distress post-operatorio e Distress pre-operatorio) e le 2 variabili statisticamente significative estratte dall'analisi fattoriale

Variabili indipendenti	Coefficienti NON standardizzati		Coefficienti standardizzati	T	Sig.
	B	Std. Error	BETA		
Costante	-1,21	0,22		-5,44	0,000
Componente fattoriale 1	-0,45	0,22	-0,17	-2,01	0,046
Componente fattoriale 3	0,59	0,22	0,22	2,65	0,009

- a. Variabile dipendente: DeltaDistress (Differenza tra Distress post-operatorio e Distress pre-operatorio)

5. Discussione e conclusioni

Nel presente lavoro, l'obiettivo fondamentale della complessiva domanda di ricerca postasi era quello di esaminare, sulla base della relazione riscontrata in letteratura tra coping e distress (Dunkel-Schetter et al., 1999), se e quanto lo stress psicologico delle

donne affette da neoplasia maligna della mammella, rilevato nel pre e post-intervento chirurgico, potesse statisticamente variare in funzione delle strategie di coping adottate nella prima fase di rilevazione. In altri termini, il quesito insorto a monte era: esiste una relazione tra le variabili “strategie adattive e non” e andamento del distress? È possibile affermare che le prime influenzino quest’ultima? E se sì, è possibile prevedere il decremento/accrecimento della quota di stress psicologico conoscendo la quantità e la tipologia di modalità messe in atto nel momento antecedente la successiva nuova valutazione dei suoi livelli?

Confermando la prima ipotesi, il distress medio era significativamente e statisticamente diminuito nella fase post-chirurgica, con un passaggio, in termini di mera frequenza, dal 40 al 23% di donne con distress psicologico clinicamente elevato. Tale riscontro è spiegabile in virtù del fatto che, mentre la fase pre-operatoria risente delle reazioni di shock e di attribuzione di connotati catastrofici conferiti subito dopo la diagnosi e, dunque, i livelli di stress psicologico risultano essere più elevati perché legati, in questo stadio, all’inquietudine principale: la morte (Worster & Holmes, 2009; Piot-Ziegler et al., 2010); la fase post-operatoria, nonostante possa essere intimamente vissuta come un’invasione, grazie all’equilibrio raggiunto tra la minaccia del cambiamento posta dalla precedente fase e le risorse di coping messe in campo per affrontarla, assume le fattezze di uno stadio di elaborazione. Si tratta, quest’ultimo, di un periodo che può essere paragonato alla “ricostruzione” in cui le donne, tentando di ripristinare nuovi obiettivi evolutivi lasciandosi alle spalle l’idea angosciante di morte, potranno indicare minori livelli di stress psicologico conseguenti alla percezione di aver superato l’ostacolo più importante dell’intero percorso diagnostico: liberarsi dal senso angosciante del “male oscuro”. In termini metaforici e psicodinamici, è come se l’operazione fosse consistita nell’asportazione, accanto alla ghiandola mammaria, della disperazione precedentemente incarnata nel corpo.

Il periodo che precede l’ospedalizzazione in attesa dell’intervento chirurgico è costellato da stati di tensione emotiva, irritabilità e incertezza: nella psiche delle pazienti, allo spettro della morte e alla percezione del limite, si aggiungono fantasie strazianti legati al timore del dolore post-operatorio, dalle complicanze che potrebbero rendere necessario un nuovo intervento, dalle preoccupazioni per il successo della chirurgia, dai rischi legati all’anestesia, dall’inaccettabilità della perdita dell’organo, dall’intrusività del pensiero legato all’immagine dell’intervento che ci si accinge a sostenere. L’aver acquisito la consapevolezza di dover continuare a lottare per sopravvivere permette, nella fase post-operatoria, di spostare il focus dell’attenzione su comportamenti strettamente connessi al ripristino della salute al fine di sentirsi completamente liberi dalla prigionia del cancro.

Parzialmente sono state, invece, confermate le ipotesi 2-3. Si era assunto che, mentre le strategie disadattive di coping (avoidant coping) sarebbero state in relazione statisticamente positiva con l'aumento di stress psicologico, predicendone l'accrescimento (H2), le modalità di approach coping, adattive, sarebbero state negativamente associate alla diminuzione dello stesso, presagendone l'abbassamento (H3) (Bussell & Naus, 2010; Groarke A, et al., 2013; Hodges K. & Winstanley S., 2012; Karademas E.C. et al., 2007). Entrambe, in ogni caso, si assumevano essere variabili in grado di predire la variazione del distress psicologico riscontrato nella fase post-operatoria.

Rispetto alla terza ipotesi, sono state le pazienti che hanno utilizzato nella fase pre-operatoria un considerevole numero di risorse di approach coping adattive, riassumibili nella prima delle sei componenti estrapolate dall'analisi fattoriale, ad essere andate incontro ad una diminuzione del distress nella fase post-chirurgica perché risultate essere, tali variabili indipendenti, in correlazione negativamente significativa con la variabile dipendente DeltaDistress. Erano le donne che hanno agito affidandosi agli altri chiedendo aiuto, consigli e supporto emotivo, che sono state in grado di rivalutare l'evento stressante in un'ottica positiva e di convivere con le avversità legate al periodo pre-operatorio e che hanno elaborato strategie cognitive idonee a migliorarlo, ad aver mostrato una riduzione dello stress psicologico. Ergo, le strategie di coping adattive possono essere variabili in grado di predire e spiegare la diminuzione riscontrata nella fase successiva all'intervento (Stanton, Danoff-Burg & Huggins, 2002).

Le analisi dei dati hanno, invece smentito la seconda ipotesi, risultando non in linea con i risultati riscontrati in letteratura, verosimilmente in quanto la seconda componente individuata dall'analisi fattoriale, era sostanzialmente legata a soltanto due delle strategie disadattive quali la negazione e il disimpegno comportamentale, un numero relativamente basso per poter raggiungere un qualche livello di significatività. Non è stata così avvalorata l'evidenza secondo cui le strategie di evitamento (avoidance coping) possano fungere da predittori di aumento del distress (Ben-Zur, Gilbar, & Lev, 2001).

Nonostante l'ipotesi sia stata rovesciata, un'evidenza statisticamente significativa che corrobora predizioni simili nell'esito ottenuto, seppur con tipologie di variabili indipendenti differenti, è emersa. L'analisi di regressione lineare multivariata ha dimostrato come fosse un terzo fattore predittivo, positivamente associato a strategie come umorismo, ristrutturazione positiva e autoaccusa ad essere positivamente correlato con le variazioni del distress, presagendone l'aumento. Si tratta di singole strategie autonome in termini categoriali in quanto afferiscono a classificazioni

diversificate tra loro, la prima spuria, la seconda di valenza adattiva e la terza disfunzionale. Malgrado quindi non si tratti di strategie globalmente etichettate come disfunzionali, tale esito spiegherebbe in ogni caso come sia la riduzione di un cospicuo numero di strategie adattive adoperate, soltanto 1 di 6, a spiegare l'aumento dello stress post-operatorio.

L'effetto riscontrato, dunque, potrebbe derivare dal fatto che una parte del campione abbia utilizzato singole modalità incongruenti dal punto di vista categoriale, perdendo di flessibilità garantita dalla numerosa messa in atto di modalità funzionali e aumentando la rigidità comportamentale. Il senso di incoerenza che lega tali variabili potrebbe essere il rispecchiamento dello sbilanciamento tra la sensazione di possedere un senso di controllo tale da attribuire alla situazione stressante una connotazione positiva e la percezione di essere essi stessi la causa dell'evento. Il tutto modulato dalla tendenza a cogliere gli elementi di ilarità dell'esperienza stressante, strategia sintomatica dell'impotenza percepita. Il sussistere di tale squilibrio situazionale si ripercuote nel grado di incongruenza manifestato nel modo di aderire alle modalità di coping. Tale incoerenza potrebbe essere spiegata come il riflesso della presenza di elevati livelli di incertezza legati sia a tratti disposizionali sia ai numerosi timori annessi all'imminente effettuazione dell'intervento chirurgico che, riversandosi sulle modalità di affrontare la situazione, hanno innescato un aumento dello stress successivo.

I risultati ottenuti nel presente lavoro scientifico, analizzati nella loro globalità, presentano importanti implicazioni cliniche nell'assessment delle donne con carcinoma mammario: le modalità adattive di far fronte all'evento "pre-operazione" valutato come stressante, sembrano essere cruciali nello spiegare la variazione del distress psicologico riscontrata nella fase conseguente alla rimozione mammaria, innescandone l'abbassamento. Pertanto, nella scelta di implementazione di interventi efficaci finalizzati a sostenere il processo di cura e sostegno, tali strategie dovrebbero essere considerate come uno dei fattori su cui intervenire al fine di disporre di un articolato repertorio di abilità di coping che possano contribuire a migliorare lo stato di benessere e fronteggiare con successo gli eventi che si susseguono nel corso dell'intero iter diagnostico.

Offrire un'adeguata assistenza alle pazienti oncologiche implica anche rilevare precocemente la presenza di potenziali disfunzionalità di carattere psicologico. La possibilità di valutare il distress delle pazienti con carcinoma della mammella durante il percorso diagnostico rappresenta l'opportunità di rilevare tempestivamente il disagio, così come dimostrato nel suddetto lavoro di ricerca. La fase pre-operatoria, infatti, risulta essere particolarmente adatta per iniziare a valutare e, di

conseguenza, prevenire lo stress psicologico, identificando le pazienti a rischio al fine di proporre loro un modello di intervento adeguato.

Alla luce dei dati ottenuti si ritiene opportuno includere uno screening di routine per il distress e far riferimento a programmi psicologici e sociali che incrementino reazioni funzionali e rispondano alle specifiche esigenze delle pazienti perché non va dimenticato che ciascuna donna rappresenta un parametro essenziale perché caratterizzata da una storia personale, dalle singolari abilità, conoscenze, aspettative. Ogni paziente ha il suo modo di sentire, pensare e comportarsi, saper identificare le peculiarità individuali e i relativi sistemi di riferimento, più o meno supportanti, è il primo passo per discernere il curare dal “prendersi cura”.

Riferimenti bibliografici

- Arraras, J. I., Wright, S. J., Jusue, G., Tejedor, M., & Calvo, J. I. (2002). Coping style, locus of control, psychological distress and pain-related behaviours in cancer and other diseases. *Psychology, Health & Medicine*, 7(2), 181-187.
- Aspinwall, L. G., & Taylor, S. E. (1997). A stitch in time: Self-regulation and proactive coping. *Psychological bulletin*, 121(3), 417.
- Ben-Zur, H., Gilbar, O., & Lev, S. (2001). Coping with breast cancer: Patient, spouse, and dyad models. *Psychosomatic Medicine*, 63(1), 32-39.
- Bower, J. E. (2008). Behavioral symptoms in breast cancer patients and survivors: fatigue, insomnia, depression, and cognitive disturbance. *Journal of clinical oncology: official journal of the American Society of Clinical Oncology*, 26(5), 768.
- Bussell, V. A., & Naus, M. J. (2010). A longitudinal investigation of coping and posttraumatic growth in breast cancer survivors. *Journal of psychosocial oncology*, 28(1), 61-78.
- Carlson, L. E., & Bultz, B. D. (2003). Benefits of psychosocial oncology care: Improved quality of life and medical cost offset. *Health and Quality of Life Outcomes*, 1(1), 1-9.
- Carver, C. S. (1997). You want to measure coping but your protocol's too long: Consider the brief cope. *International journal of behavioral medicine*, 4(1), 92.
- Carver, C. S., Scheier, M. F., & Weintraub, J. K. (1989). Assessing coping strategies: a theoretically based approach. *Journal of personality and social psychology*, 56(2), 267.
- Conley, C. C., Bishop, B. T., & Andersen, B. L. (2016, September). Emotions and emotion regulation in breast cancer survivorship. In *Healthcare* (Vol. 4, No. 3, p. 56). Multidisciplinary Digital Publishing Institute.

- Conti, L. (1999). *Repertorio delle scale di valutazione in psichiatria-Tomo I*, SEE Firenze
- Dunkel-Schetter, C., Feinstein, L. G., Taylor, S. E., & Falke, R. L. (1999). Patterns of coping with cancer. In: R. M. Suinn & G. R. VandenBos (eds.), *Cancer patients and their families: Readings on disease course, coping, and psychological interventions*. American Psychological Association, 35–51.
- Endler, N. S., & Parker, J. D. (1990). Multidimensional assessment of coping: A critical evaluation. *Journal of personality and social psychology*, 58(5), 844.
- Folkman, S., & Lazarus, R. S. (1984). *Stress, appraisal, and coping* (pp. 150-153). New York: Springer Publishing Company.
- Gil, F., Grassi, L., Travado, L., Tomamichel, M., Gonzalez, J. R., & Sepos Group. (2005). Use of distress and depression thermometers to measure psychosocial morbidity among southern European cancer patients. *Supportive care in cancer*, 13(8), 600-606.
- Gilbar, O., Or-Han, K., & Plivazky, N. (2005). Mental adjustment, coping strategies, and psychological distress among end-stage renal disease patients. *Journal of Psychosomatic Research*, 58(6), 471-476.
- Glinder, J. G., Beckjord, E., Kaiser, C. R., & Compas, B. E. (2007). Psychological adjustment to breast cancer: Automatic and controlled responses to stress. *Psychology and Health*, 22(3), 337-359.
- Groarke, A., Curtis, R., & Kerin, M. (2013). Cognitive-behavioural stress management enhances adjustment in women with breast cancer. *British journal of health psychology*, 18(3), 623-641.
- Hodges, K., & Winstanley, S. (2012). Effects of optimism, social support, fighting spirit, cancer worry and internal health locus of control on positive affect in cancer survivors: a path analysis. *Stress and Health*, 28(5), 408-415.
- Holland, J. C., & Alici, Y. (2010). Management of distress in cancer patients. *The journal of supportive oncology*, 8(1), 4.
- Karademas, E. C., Karvelis, S., & Argyropoulou, K. (2007). Stress-related predictors of optimism in breast cancer survivors. *Stress and Health: Journal of the International Society for the Investigation of Stress*, 23(3), 161-168.
- Mehnert, A., Scherwath, A., Schirmer, L., Schleimer, B., Petersen, C., Schulz-Kindermann, F., ... & Koch, U. (2007). The association between neuropsychological impairment, self-perceived cognitive deficits, fatigue and health related quality of life in breast cancer survivors following standard adjuvant versus high-dose chemotherapy. *Patient Education and Counseling*, 66(1), 108-118.
- Meijer, A., Roseman, M., Delisle, V. C., Milette, K., Levis, B., Syamchandra, A., ... & Thombs, B. D. (2013). Effects of screening for psychological distress on patient outcomes in cancer: a systematic review. *Journal of psychosomatic research*, 75(1), 1-17.

- Montgomery, C., Pocock, M., Titley, K., & Lloyd, K. (2003). Predicting psychological distress in patients with leukaemia and lymphoma. *Journal of psychosomatic research*, 54(4), 289-292.
- Oh, S., Miyamoto, H., Yamazaki, A., Fukai, R., Shiomi, K., Sonobe, S., ... & Sakao, Y. (2007). Prospective analysis of depression and psychological distress before and after surgical resection of lung cancer. *General thoracic and cardiovascular surgery*, 55(3), 119-124.
- Oswiecki, D. M., & Compas, B. E. (1999). A prospective study of coping, perceived control, and psychological adaptation to breast cancer. *Cognitive therapy and research*, 23(2), 169-180.
- Piot-Ziegler, C., Sassi, M. L., Raffoul, W., & Delaloye, J. F. (2010). Mastectomy, body deconstruction, and impact on identity: a qualitative study. *British Journal of Health Psychology*, 15(3), 479-510.
- Pugliese, P., Falcicchio, C. (2015). Valutazione dei bisogni di salute psicosociale e di supporto. In: Antonuzzo A, et al. *Manuale di cure di supporto in oncologia*. Roma: Universo.
- Stanton, A. L., Danoff-burg, S., & Huggins, M. E. (2002). The first year after breast cancer diagnosis: hope and coping strategies as predictors of adjustment. *Psycho-Oncology: Journal of the Psychological, Social and Behavioral Dimensions of Cancer*, 11(2), 93-102.
- Suls, J., & Fletcher, B. (1985). The relative efficacy of avoidant and nonavoidant coping strategies: a meta-analysis. *Health psychology*, 4(3), 249.
- Veronesi, U., Cascinelli, N., Mariani, L., Greco, M., Saccozzi, R., Luini, A., ... & Marubini, E. (2002). Twenty-year follow-up of a randomized study comparing breast-conserving surgery with radical mastectomy for early breast cancer. *New England Journal of Medicine*, 347(16), 1227-1232.
- Vos, P. J., Garssen, B., Visser, A. P., Duivenvoorden, H. J., & de Haes, H. C. (2004). Early stage breast cancer: Explaining level of psychosocial adjustment using structural equation modeling. *Journal of Behavioral Medicine*, 27(6), 557-580.
- Worster, B., & Holmes, S. (2009). A phenomenological study of the postoperative experiences of patients undergoing surgery for colorectal cancer. *European Journal of Oncology Nursing*, 13(5), 315-322.

Valutazione dei servizi web offerti dall'Università di Foggia ai tempi del COVID-19*

Laura Antonucci, Corrado Crocetta, Yana Kostiuk

Università degli Studi di Foggia

Riassunto: Nel presente elaborato vengono presentati i risultati dell'indagine sulla customer satisfaction dei servizi web offerti dall'Università di Foggia durante la pandemia da COVID-19. L'universo di riferimento dell'indagine è rappresentato dagli studenti iscritti all'Università di Foggia nell'anno accademico 2019-20. Attraverso l'applicazione dei modelli PLS-PM vengono individuate le variabili di maggiore criticità a cui dedicare maggiore attenzione al fine di soddisfare al meglio le esigenze degli studenti.

Keywords: servizi web; customer satisfaction; modelli ad equazioni strutturali; PLS-PM.

1. Introduzione

La qualità dei siti web diventa sempre più determinante soprattutto nel periodo di emergenza da Covid-19, quando il web diventa indispensabile e, in alcuni casi, vitale per il commercio, comunicazione, intrattenimento o qualsiasi altro tipo di attività facente parte della routine giornaliera delle persone. L'emergenza ha colpito anche le università che hanno dovuto reagire e adattarsi, passando alla modalità online in tutti gli ambiti. In questo contesto la qualità dei servizi web non solo influenza la qualità dell'istruzione, ma rappresenta allo stesso tempo sia una minaccia che un'opportunità per la competitività delle università.

* Il presente articolo è frutto del lavoro congiunto degli autori, tuttavia la stesura finale dei paragrafi 3. , 3.1 e 3.2 è da attribuirsi a L. Antonucci, dei paragrafi 1 e 2 a C. Crocetta, e a Y. Kostiuk quella dei paragrafi 3.3., 3.4 e 4.

Alla luce di queste considerazioni, l'obiettivo principale di questo elaborato è la valutazione del livello di soddisfazione degli studenti rispetto ai servizi web offerti dall'Università di Foggia ai tempi del Covid-19.

In particolare, la ricerca è incentrata su tre siti:

- Sito web ufficiale (unifg.it);
- Segreteria online (esse3.unifg.it);
- Portale dei servizi e-learning (elearning.unifg.it).

2. Materiali e metodi

Al fine di valutare il livello di soddisfazione degli studenti rispetto ai servizi web dell'Università di Foggia è stato costruito un questionario online, inoltrato a tutti gli studenti attraverso la posta elettronica istituzionale. Il questionario è composto da 25 domande, suddivise in cinque sezioni:

- Profilo studente;
- Usabilità;
- Contenuti;
- Grafica;
- Manutenzione e giudizio complessivo.

La somministrazione è avvenuta on-line nel periodo dal 24/06/2020 al 14/07/2020. In totale sono pervenute 1.063 risposte valide.

I dati ottenuti sono stati utilizzati per lo sviluppo di un modello di tipo PLS-PM al fine di stimare le relazioni esistenti tra le variabili che definiscono la qualità di un sito e di individuare quella avente maggior impatto sul livello di soddisfazione degli studenti.

I risultati del presente elaborato sono stati ottenuti riportando all'universo i dati campionari, ottenuti dal questionario, attraverso un coefficiente di espansione dato da:

$$w_{ijk} = N_{ijk}/n_{ijk}$$

dove w_{ijk} è il coefficiente di espansione, N_{ijk} è la numerosità della popolazione degli studenti iscritti all'Università di Foggia nell'anno accademico 2019-20 suddivisi per *Dipartimento*, *Corso di studio* e *Anno di corso* e n_{ijk} è il campione realizzato degli studenti, stratificato per *Dipartimento*, *Corso di studio* e *Anno di corso*.

Nelle tabelle di seguito vengono riportati i coefficienti di espansione ottenuti.

Tabella 1. *Coefficienti di espansione per le lauree triennali*

	1 anno	2 anno	3 anno	FC
Economia	10,9	10,7	8,9	11,3
Giurisprudenza	17,2	22,2	8,6	22,2
Medicina CS	9,8	18,7	11,5	26,6
Scienze Agrarie	11,6	7,1	8,27	18,5
Scienze MC	10,8	11,4	6,8	61,0
Studi Umanistici	9,3	5,7	5,0	12,2

Tabella 2. *Coefficienti di espansione per le lauree magistrali*

	1 anno	2 anno	FC
Economia	6,1	8,4	5,9
Giurisprudenza	10,7	-	-
Scienze Agrarie	7,3	5,6	18,7
Studi Umanistici	7,1	6,7	16,4

Tabella 3. *Coefficienti di espansione per le lauree magistrali a ciclo unico*

	1 anno	2 anno	3 anno	4 anno	5 anno	6 anno	FC
Giurisprudenza	25,8	17,7	12,2	15,8	16,3	-	28,1
Scienze MC	9,1	4,6	3,7	7,6	5,3	9,1	7,7
Medicina CS	5,6	3,6	8,5	5,3	7,3	14,7	5,0

Il coefficiente più critico (con valore elevato) riguarda gli studenti fuori corso delle lauree triennali del Dipartimento di Scienze Mediche e Chirurgiche (pari a 61).

3. Risultati

Gli studenti hanno assegnato un voto da 1 a 5 per esprimere il giudizio complessivo sui tre siti oggetto della ricerca.

Per il sito web ufficiale, il voto più basso è pari a 3.4, dato dagli studenti dei Dipartimenti di Medicina Clinica e Sperimentale e di Studi Umanistici. Il voto più basso per la Segreteria on-line, pari a 3.5, è stato assegnato dagli studenti di Studi Umanistici. Infine, la valutazione del Portale dei servizi e-learning risulta più positiva rispetto agli altri due siti. Infatti, il punteggio più basso è pari a 3.7, assegnato dagli studenti di Giurisprudenza (Tab.4).

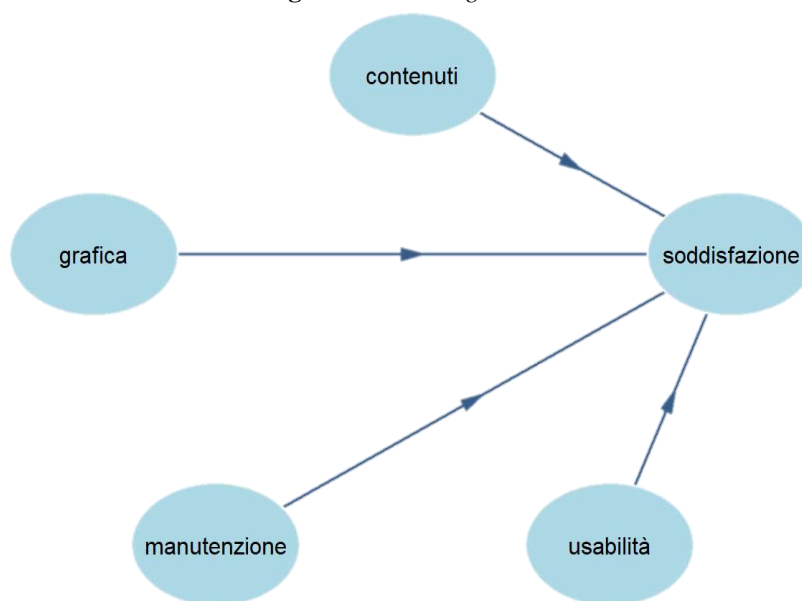
Tabella 4. Giudizi medi in base alle risposte espressi dai rispondenti sui 3 servizi analizzati, per Dipartimento

Servizio analizzato	Economia	Giurisprud.	Medicina CS	Scienze Agrarie	Scienze MC	Studi Umanistici
Sito web ufficiale	3,8	3,5	3,4	3,7	3,5	3,4
Segreteria online	3,8	3,6	3,6	3,9	3,6	3,5
Portale dei servizi e-learning	3,8	3,7	3,8	3,9	3,9	3,8

Per l'analisi dei dati ottenuti dal questionario è stato scelto l'approccio di soft modeling PLS-PM, ampiamente applicato nell'ambito della Customer satisfaction. L'obiettivo di quest'analisi è quello di poter definire la Customer satisfaction (in questo caso la soddisfazione degli studenti rispetto ai servizi web offerti dall'università) come una variabile influenzata da alcune variabili latenti, stimando quale parte della sua variabilità il modello riesce a spiegare.

Il modello è stato creato collegando le 4 variabili latenti esogene: "Usabilità", "Contenuti", "Grafica" e "Manutenzione" con la variabile endogena "Soddisfazione" (Fig. 1).

Figura 1. Path diagram



Nella tabella di seguito vengono riportate le variabili manifeste considerate per ogni variabile latente.

Tabella 5. *Variabili manifeste e relative variabili latenti*

Etichetta	Variabile Manifesta	Variabile Latente
<i>acc</i>	Facilità di accesso	Usabilità
<i>nav</i>	Facilità di navigazione	
<i>mob</i>	Mobile-friendly	
<i>brow</i>	Accessibilità qualsiasi browser	
<i>men_org</i>	Organizzazione menù navigazione	
<i>temp</i>	Tempo caricamento pagine	
<i>cont_necess</i>	Presenza di contenuti necessari	Contenuti
<i>cont_completi</i>	Completezza di contenuti	
<i>cont_aggiornati</i>	Aggiornamento contenuti	
<i>avvisi</i>	Reperibilità degli avvisi	
<i>grafica</i>	Grafica	Grafica
<i>caratt_leg</i>	Leggibilità testo	
<i>interr_accett</i>	Durata interruzioni	Manutenzione
<i>interr_inf</i>	Adeguatezza informativa	
<i>compless</i>	Giudizio complessivo	Soddisfazione

3.1 *Sito web ufficiale*

Il modello di misura scelto è quello di tipo riflessivo. Tale scelta è dovuta alla presenza di multicollinearità tra gli indicatori. La scelta del modello di tipo riflessivo è avvalorata dal fatto che gli indici quali l'Alpha di Cronbach, il Rho di Dillon-Goldstein e il primo autovalore confermano l'unidimensionalità dei blocchi. Infatti, i primi due indicatori risultano superiori a 0,7, mentre il terzo è maggiore di 1 (Tab.6).

Tabella 6. *Consistenza interna dei blocchi (Sito web ufficiale)*

	Modo	N.VM	C.alpha	DG.rho	1°A.v.	2°A.v.
Usabilità	A	5	0,92	0,94	3,76	0,47
Contenuti	A	4	0,90	0,93	3,11	0,44
Grafica	A	2	0,84	0,93	1,73	0,27
Manutenzione	A	2	0,78	0,90	1,64	0,36
Soddisfazione	A	1	1	1	1	0

Per valutare la qualità del modello sono stati considerati i loadings e le comunalità relative a ogni indicatore.

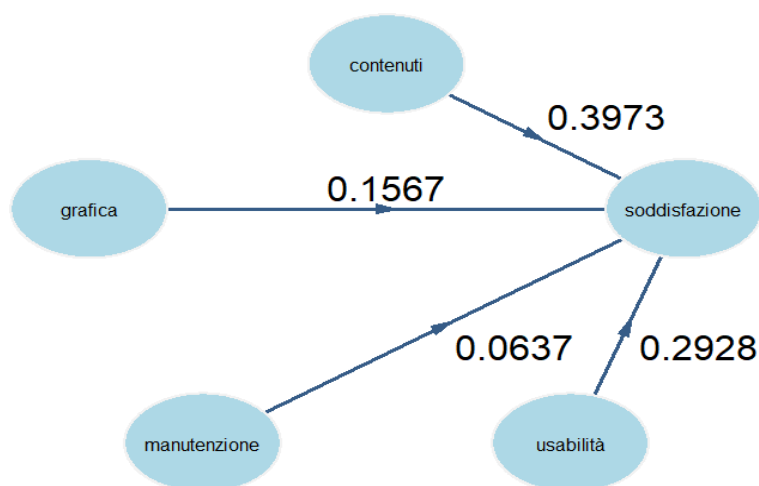
Tabella 7. Loadings e comunalità (Sito web ufficiale)

VL	VM	Loading	Comunalità	AVE
Usabilità	<i>acc</i>	0,88	0,79	0,75
	<i>mob</i>	0,83	0,69	
	<i>temp</i>	0,81	0,66	
	<i>nav</i>	0,91	0,82	
	<i>men_org</i>	0,89	0,79	
Contenuti	<i>cont_necess</i>	0,91	0,83	0,78
	<i>cont_completi</i>	0,92	0,84	
	<i>cont_aggiornati</i>	0,87	0,76	
	<i>avvisi</i>	0,82	0,68	
Grafica	<i>grafica</i>	0,94	0,88	0,86
	<i>caratt_leg</i>	0,92	0,85	
Manutenzione	<i>interr_accett</i>	0,89	0,79	0,82
	<i>interr_inf</i>	0,92	0,85	
Soddisfazione	<i>compless</i>	1	1	1

Come si evince dalla Tab.7, tutti i loadings e le comunalità hanno i valori molto alti, per cui la validità convergente del modello è confermata. Inoltre, analizzando le comunalità medie (AVE), possiamo notare che la variabile *Grafica* è quella che spiega nel modo migliore la variabilità dei propri indicatori (in media l'86%).

Nella figura 2 è riportata una rappresentazione grafica del modello interno.

Figura 2. Path diagram (Sito web ufficiale)



Dalla Fig. 2 si evince che la VL *Contenuti* è quella avente il maggior impatto (0.397), seguita dalla VL *Usabilità* (0.29); mentre la variabile *Manutenzione* ha l'impatto residuale (0.06). Per verificare la significatività dei coefficienti è stata utilizzata la procedura di bootstrap. Gli intervalli di confidenza stimati non includevano 0, per cui è stata confermata la significatività dei path coefficients al 5%.

Riguardo al modello esterno, i valori da considerare sono i pesi esterni che permettono di valutare l'importanza di ogni singolo indicatore nella formazione della corrispondente variabile latente. Come si nota dalla Tab.8, tutti gli indicatori contribuiscono quasi equamente alla determinazione della corrispondente variabile latente.

Tabella 8. *Pesi esterni (Sito web ufficiale)*

VL	VM	Peso esterno
Usabilità	<i>acc</i>	0,227
	<i>mob</i>	0,205
	<i>temp</i>	0,205
	<i>nav</i>	0,254
	<i>men_org</i>	0,259
Contenuti	<i>cont_necess</i>	0,305
	<i>cont_completi</i>	0,312
	<i>cont_aggiornati</i>	0,259
	<i>avvisi</i>	0,255
Grafica	<i>grafica</i>	0,571
	<i>caratt_leg</i>	0,506
Manutenzione	<i>interr_accett</i>	0,499
	<i>interr_inf</i>	0,603

Infine, bisogna analizzare i punteggi delle variabili latenti, che rappresentano una stima del giudizio dei rispondenti sui vari costrutti latenti.

Tabella 9. *Punteggi delle variabili latenti (Sito web ufficiale)*

	Minimo	Massimo	Media	Dev.st.
Usabilità	1	5	3,66	1,09
Contenuti	1	5	3,35	1,13
Grafica	1	5	3,76	1,10
Manutenzione	1	5	3,31	1,13
Soddisfazione	1	5	3,54	1,11

Come si nota dalla Tab.9, la variabile *Grafica* ha il punteggio più alto, pari a 3.76; mentre la variabile avente il punteggio più basso è *Manutenzione* (3.31). La variabile di interesse principale, ovvero *Soddisfazione*, ha un punteggio stimato pari a 3.54. Il coefficiente R^2 , risulta essere pari a 0.69, dunque il modello riesce a spiegare il 69% della variabilità della VL *Soddisfazione*. Infine, per quanto riguarda il modello globale, l'indice GoF, pari a 0.74, risulta superiore al valore accettabile di 0.7.

3.2 *Segreteria online*

Come si nota dalla tabella 10, tutti gli indicatori superano la soglia prefissata, per cui l'unidimensionalità dei blocchi è confermata.

Tabella 10. *Consistenza interna dei blocchi (Segreteria online)*

	Modo	N.VM	C.alpha	DG.rho	1°A.v.	2°A.v.
Usabilità	A	5	0,93	0,95	3,90	0,37
Contenuti	A	4	0,92	0,94	3,20	0,40
Grafica	A	2	0,85	0,93	1,74	0,26
Manutenzione	A	2	0,79	0,90	1,65	0,35
Soddisfazione	A	1	1	1	1	0

Anche la validità convergente del modello esterno è confermata, in quanto i loadings superano la soglia di 0.7 e, di conseguenza, anche tutte le comunaltà hanno i valori accettabili, ovvero superiori a 0.5 (Tab.11).

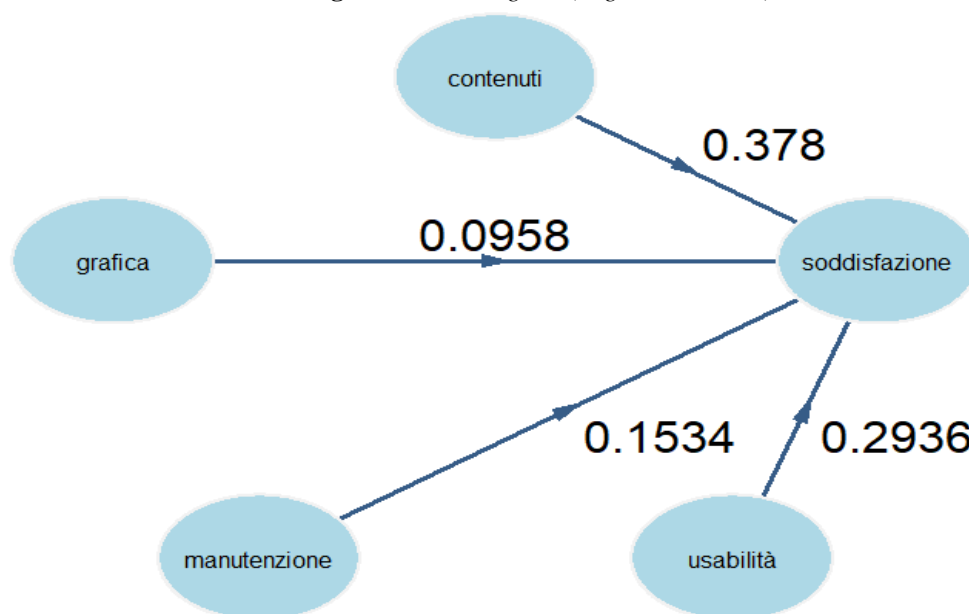
Tabella 11. *Loadings e comunaltà (Sito web ufficiale)*

VL	VM	Loading	Comunalità	AVE
Usabilità	<i>acc</i>	0,90	0,80	0,78
	<i>mob</i>	0,85	0,73	
	<i>temp</i>	0,84	0,71	
	<i>nav</i>	0,91	0,83	
	<i>men_org</i>	0,91	0,83	
Contenuti	<i>cont_necess</i>	0,91	0,84	0,80
	<i>cont_completi</i>	0,93	0,86	
	<i>cont_aggiornati</i>	0,90	0,80	
	<i>avvisi</i>	0,83	0,70	
Grafica	<i>grafica</i>	0,94	0,88	0,87
	<i>caratt_leg</i>	0,93	0,86	
Manutenzione	<i>interr_accett</i>	0,91	0,82	0,83
	<i>interr_inf</i>	0,91	0,83	
Soddisfazione	<i>compless</i>	1	1	1

Analizzando la comunalità media, si nota che tutte le variabili, tranne *Usabilità*, spiegano più del 80% della variabilità dei loro indicatori.

Analizzando il modello interno, rappresentato dal path diagram (Fig.3), possiamo notare che la variabile *Contenuti* è quella avente il maggior impatto sulla soddisfazione (0.38), seguita dalla variabile *Usabilità* (0.29); mentre la variabile *Grafica* risulta la meno importante (0.096).

Figura 3. Path diagram (*Segreteria online*)



Per quanto riguarda il modello esterno, come da Tab.12, tutti gli indicatori contribuiscono quasi ugualmente alla costruzione della corrispondente variabile latente. Gli intervalli di confidenza, stimati con l'ausilio del Bootstrap, non comprendono 0, per cui possiamo confermare la significatività dei pesi esterni al 5%.

Tabella 12. Pesi esterni (Segreteria online)

VL	VM	Peso esterno
Usabilità	<i>acc</i>	0,227
	<i>mob</i>	0,217
	<i>temp</i>	0,209
	<i>nav</i>	0,234
	<i>men_org</i>	0,244
Contenuti	<i>cont_necess</i>	0,293
	<i>cont_completi</i>	0,297
	<i>cont_aggiornati</i>	0,275
	<i>avvisi</i>	0,252
Grafica	<i>grafica</i>	0,555
	<i>caratt_leg</i>	0,519
Manutenzione	<i>interr_accett</i>	0,539
	<i>interr_inf</i>	0,560

Tabella 13. Punteggi delle variabili latenti (Segreteria online)

	Minimo	Massimo	Media	Dev.st.
Usabilità	1	5	3,77	1,06
Contenuti	1	5	3,58	1,12
Grafica	1	5	3,81	1,09
Manutenzione	1	5	3,20	1,14
Soddisfazione	1	5	3,54	1,11

Infine, *analizziamo* i punteggi stimati (scores) delle variabili latenti.

Come si nota dalla Tab.13, la variabile latente con lo score maggiore è ancora una volta *Grafica* (3.81 in media), mentre la variabile con il minor punteggio è *Manutenzione* (3.2). La soddisfazione presenta un punteggio medio pari a 3.66, indicando un livello più alto rispetto al sito valutato precedentemente, ma comunque migliorabile.

Il coefficiente R^2 risulta essere pari a 0.69, equivale a dire che il 69% della variabilità della variabile latente Soddisfazione è spiegato dal modello.

Infine, la capacità predittiva del modello, data dall'indice GoF, è pari a 0.74.

3.3 *Portale dei servizi e-learning*

Analogamente ai due siti analizzati prima, l'unidimensionalità dei blocchi è confermata da tutti gli indici. Infatti, l'Alpha di Cronbach e il Rho di Dillon-

Goldstein risultano superiori a 0.7 e il primo autovalore è sempre maggiore di 1 (Tab.14).

Tabella 14. *Consistenza interna dei blocchi (Portale dei servizi e-learning)*

	Modo	N.VM	C.alpha	DG.rho	1°A.v.	2°A.v.
Usabilità	A	5	0,92	0,94	3,85	0,40
Contenuti	A	4	0,92	0,94	3,23	0,44
Grafica	A	2	0,87	0,94	1,76	0,24
Manutenzione	A	2	0,80	0,91	1,67	0,33
Soddisfazione	A	1	1	1	1	0

I loadings e le comunalità indicati nella tabella sotto ci permettono di confermare la validità convergente del modello. Infatti, tutti i loadings sono maggiori di 0.7 e tutte le comunalità risultano superiori a 0.5.

Tabella 15. *Loadings e comunalità (Portale dei servizi e-learning)*

VL	VM	Loading	Comunalità	AVE
Usabilità	<i>acc</i>	0,90	0,81	0,77
	<i>mob</i>	0,80	0,66	
	<i>temp</i>	0,83	0,70	
	<i>nav</i>	0,93	0,86	
	<i>men_org</i>	0,91	0,82	
Contenuti	<i>cont_necess</i>	0,93	0,86	0,81
	<i>cont_completi</i>	0,94	0,88	
	<i>cont_aggiornati</i>	0,92	0,84	
	<i>avvisi</i>	0,8	0,65	
Grafica	<i>grafica</i>	0,94	0,89	0,87
	<i>caratt_leg</i>	0,94	0,88	
Manutenzione	<i>interr_accett</i>	0,92	0,85	0,84
	<i>interr_inf</i>	0,91	0,82	
Soddisfazione	<i>compless</i>	1	1	1

Come si evince dalla Tab.15, la variabile *Grafica* spiega gran parte della variabilità dei corrispondenti indicatori (87% in media), mentre l'Usabilità, così come nei due casi precedenti, non supera l'80%.

Per comprendere l'importanza dei vari indicatori analizziamo i loro pesi esterni.

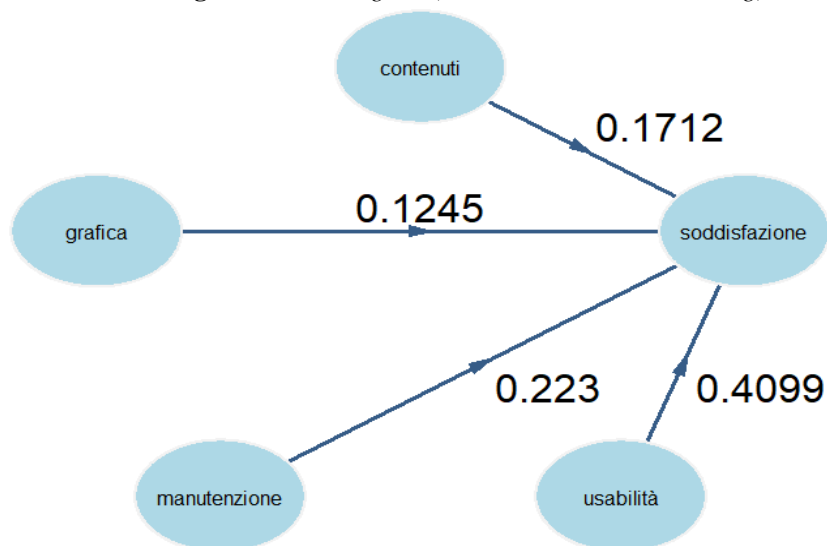
Tabella 16. Pesi esterni (Portale dei servizi e-learning)

VL	VM	Peso esterno
Usabilità	<i>acc</i>	0,227
	<i>mob</i>	0,217
	<i>temp</i>	0,209
	<i>nav</i>	0,234
	<i>men_org</i>	0,244
Contenuti	<i>cont_necess</i>	0,293
	<i>cont_completi</i>	0,297
	<i>cont_aggiornati</i>	0,275
	<i>avvisi</i>	0,252
Grafica	<i>grafica</i>	0,555
	<i>caratt_leg</i>	0,519
Manutenzione	<i>interr_accett</i>	0,539
	<i>interr_inf</i>	0,560

Come si nota dalla Tab.16, quasi tutti gli indicatori contribuiscono ugualmente alla costruzione dei corrispondenti costrutti latenti. L'unico indicatore per cui si osserva un peso leggermente minore rispetto agli altri indicatori dello stesso blocco è *mob* (possibilità di accedere dai vari dispositivi) con il peso pari a 0.197. La significatività dei pesi esterni è stata confermata dal Bootstrap.

Le stime del modello interno sono riportate nel path diagram (Fig.4).

Figura 4. Path diagram (Portale dei servizi e-learning)



Dalla figura di cui sopra risulta che la VL *Soddisfazione* risulta fortemente influenzata dalla variabile esogena *Usabilità* (0.409). La variabile *Manutenzione* ha un impatto moderato (0.223), mentre le variabili *Contenuti* e *Grafica* esercitano quasi lo stesso effetto sul livello della soddisfazione: 0.17 e 0.12 rispettivamente. A differenza degli altri due siti, in questo modello non ci sono variabili latenti con un impatto residuale.

Infine, possiamo analizzare i punteggi delle variabili latenti.

Tabella 17. *Punteggi delle variabili latenti (Portale dei servizi e-learning)*

	Minimo	Massimo	Media	Dev.st.
Usabilità	1	5	3,85	1,04
Contenuti	1	5	3,80	1,07
Grafica	1	5	3,95	1,04
Manutenzione	1	5	3,08	1,18
Soddisfazione	1	5	3,83	1,11

Come si evince dalla Tab.4.13, gli studenti risultano particolarmente soddisfatti della grafica (3.95 in media), mentre la variabile con il minor punteggio in assoluto è la *Manutenzione* (3.08). La soddisfazione presenta un punteggio medio pari a 3.83, indicando un livello più alto rispetto agli altri due siti.

Il coefficiente R^2 è pari a 0,69, per cui il 69% della variabilità della Soddisfazione è spiegato dalla sua dipendenza dagli altri costrutti latenti.

La capacità predittiva del modello, data dall'indice GoF, è pari a 0,74.

3.4 *Confronto tra i modelli stimati*

Un confronto dei modelli stimati per i tre siti oggetto di ricerca ci ha permesso di giungere ad una serie di conclusioni, riportate di seguito.

Per quanto riguarda la costruzione dei pesi esterni, che rappresentano l'importanza di un indicatore nella costruzione delle variabili latenti, come da Tab.18, è emerso che:

- nella costruzione della variabile Usabilità gli indicatori più importanti per tutti e tre siti sono *nav* (facilità di navigazione) e *men_org* (organizzazione del menù di navigazione);
- la completezza e l'organizzazione dei contenuti (*cont_completi*) è l'indicatore con il maggior peso per tutti i siti, nella costruzione della variabile Contenuti;
- nella formazione della variabile Grafica è l'aspetto visivo (*grafica*) ad avere la maggior importanza;

- nella costruzione della VL Manutenzione per la Segreteria online e il Portale dei servizi e-learning l'adeguatezza dell'informativa e la durata delle interruzioni risultano ugualmente importanti, mentre per il Sito web ufficiale l'adeguatezza informativa ha il peso maggiore.

Tabella 18. Pesi esterni a confronto

VL	VM	Sito web ufficiale	Segreteria online	Portale dei servizi e-learning
Usabilità	<i>acc</i>	0,227	0,227	0,236
	<i>mob</i>	0,205	0,217	0,197
	<i>temp</i>	0,205	0,209	0,225
	<i>nav</i>	0,254	0,234	0,241
	<i>men_org</i>	0,259	0,244	0,239
Contenuti	<i>cont_necess</i>	0,305	0,293	0,292
	<i>cont_completi</i>	0,312	0,297	0,301
	<i>cont_aggiornati</i>	0,259	0,275	0,275
	<i>avvisi</i>	0,255	0,252	0,243
Grafica	<i>grafica</i>	0,571	0,555	0,541
	<i>caratt_leg</i>	0,506	0,519	0,524
Manutenzione	<i>interr_accett</i>	0,499	0,539	0,546
	<i>interr_inf</i>	0,603	0,56	0,53

Per quanto riguarda, invece, le stime dei path coefficients, analizzando la Tab.19, possiamo fare le seguenti considerazioni:

- la VL esogena *Usabilità* esercita la stessa influenza sulla variabile d'interesse *Soddisfazione* per il Sito web ufficiale e per la Segreteria online, mentre risulta la più importante per il Portale dei servizi e-learning;
- per quanto riguarda la variabile *Contenuti*, essa risulta la più importante per il Sito web ufficiale e per la Segreteria online, mentre ha un impatto relativamente basso nella formazione della soddisfazione rispetto al Portale dei servizi e-learning;
- analizzando la variabile *Grafica*, si nota che per la Segreteria online è quella con il minor impatto, mentre per gli altri due siti risulta abbastanza importante;
- infine, la VL *Manutenzione* esercita un'influenza residuale sul livello di *Soddisfazione* rispetto al Sito web ufficiale, mentre risulta abbastanza importante per la Segreteria online e molto importante per il Portale dei servizi e-learning.

Tabella 19. *Path coefficients a confronto*

	Sito web ufficiale	Segreteria online	Portale dei servizi e-learning
Usabilità	0,293	0,29	0,409
Contenuti	0,397	0,377	0,171
Grafica	0,156	0,095	0,125
Manutenzione	0,063	0,15	0,223

Le considerazioni fatte risultano molto importanti al fine di individuare le variabili di maggiore criticità, in quanto permettono di elaborare delle strategie di miglioramento differenziate per ognuno di tre siti, tenendo conto dell'importanza dei vari indicatori nella formazione dei costrutti latenti.

4. Conclusioni

Il presente elaborato ha preso in esame la qualità dei siti web e le sue determinanti, sottolineando la sua importanza per la Customer satisfaction, in particolare nell'ambito universitario.

Per quanto riguarda il Sito web ufficiale, dallo studio è emerso che la variabile latente *Soddisfazione* dipende maggiormente dalla variabile *Contenuti* con il path coefficient pari a 0.4, seguita da *Usabilità* con il path coefficient pari a 0.3. La variabile *Manutenzione* ha un impatto trascurabile, pari a 0.06.

Riguardo alla Segreteria online, le variabili che esercitano maggiore influenza sono *Contenuti* (0.38) ed *Usabilità* (0.3). A differenza del Sito web ufficiale, è la variabile *Grafica* (0.09) ad avere l'impatto residuale, mentre la variabile *Manutenzione* ha un impatto moderato (0.15).

Infine, per il Portale dei servizi e-learning, le due dimensioni più importanti sono *Usabilità* (0.4) e *Manutenzione* (0.22), mentre le altre due hanno un effetto moderato: 0,17 per *Contenuti* e 0,12 per *Grafica*.

Riferimenti bibliografici

- Esposito Vinzi V.; Chin W.W.; Henseler J.; Wang H. (2010). *Handbook of Partial Least Squares: Concepts, Methods and Applications in Marketing and Related Fields*, Springer.
- Esposito Vinzi V.; Trinchera L.; Amato S. (2010). *PLS path modeling: from foundations to recent developments and open issues for model assessment and improvement*, Handbook of Partial Least Squares.

- Hair J.F., Ringle C.M.; Sarstedt M. (2011). *PLS-SEM: Indeed a silver bullet*. In: *Journal of Marketing Theory and Practice*, 19.
- Ringle C.M.; Sarstedt M., Hair J.F., Hult G. (2014). *A Primer on Partial Least Squares Structural Equation Modeling (PLS-SEM)*, Sage.
- Sanchez G. (2013). *PLS Path Modeling with R*, Trowchez Editions, Berkeley.
- Tenenhaus M.; Esposito Vinzi V.; Chatelin Y.; Lauro C. (2005), *PLS path modeling*, *Computational Statistics and Data Analysis* 48.

Analisi della soddisfazione rispetto ai servizi erogati da Sanitaservice ASL Foggia s.r.l.*

Laura Antonucci, Corrado Crocetta,
Yana Kostiuk, Massimo Russo

Università degli Studi di Foggia

Riassunto: Il presente elaborato ha come obiettivo l'analisi della customer satisfaction dei servizi erogati da Sanitaservice ASL FG s.r.l.. Lo studio consente di confrontare l'andamento del fenomeno relativo agli anni 2017,2018 e 2019 e di individuare i punti di forza e di debolezza per ogni servizio e per le singole dimensioni esaminate. L'indagine ha evidenziato che la Sanitaservice ASL FG ha sensibilmente migliorato le proprie capacità di risolvere problemi e fornire servizi con elevati standard qualitativi.

Keywords: valutazione servizi sanitari, customer satisfaction.

1. Introduzione

Nel 2019 l'azienda Sanitaservice ASL FGASL FG s.r.l. ha affidato al Dipartimento di Economia dell'Università di Foggia il compito di analizzare il livello di soddisfazione dei propri utenti. Tale indagine, giunta ormai alla terza edizione, mira a valutare, attraverso i giudizi espressi dai responsabili delle diverse strutture della ASL FG che si avvalgono dei servizi della Sanitaservice ASL FGASL FG, quale sia il livello di soddisfazione degli stessi e ad individuare eventuali criticità o aree di miglioramento.

* Il presente articolo è frutto del lavoro congiunto degli autori, tuttavia la stesura finale dei paragrafi 2.1, 3.1 e 3.6 è da attribuirsi a L. Antonucci, dei paragrafi 1 e 4 a C. Crocetta, dei paragrafi 3.3, 3.4 e 3.5 a Y. Kostiuk, e a M. Russo quella dei paragrafi 2.2 e 3.2.

2. Il piano dell'indagine e gli strumenti utilizzati

L'oggetto della ricerca e i diversi aspetti analizzati sono rimasti invariati rispetto alle precedenti edizioni, per ovvi motivi di continuità e per garantire la confrontabilità dei dati rilevati nei due anni precedenti. Tuttavia, nel 2019 sono state introdotte alcune novità utili per risolvere alcune criticità emerse nelle edizioni precedenti, ma soprattutto per tener conto della situazione di emergenza in cui si è svolta la rilevazione. L'indagine, infatti, è partita a metà febbraio 2020 ma la raccolta dei questionari ha avuto inizio a marzo 2020, in piena fase di emergenza COVID-19. Per questo motivo si è reso necessario affiancare ai questionari cartacei utilizzati nelle precedenti edizioni anche dei questionari online.

2.1 *Il questionario e le sue novità*

Anche nell'edizione 2019 della ricerca i servizi erogati da Sanitaservice ASL-FGASL FG analizzati a supporto alle attività istituzionali dell'ASL FG sono stati:

- Servizio di Ausiliariato;
- Servizio di Pulizie e Sanificazione ambientale;
- Servizio di Manutenzione degli immobili;
- Servizio di Trasporto beni economici;
- Servizio di Ausiliariato dell'UU.DD.TT.;
- Servizio Infermieristico dell'UU.DD.TT.;
- Servizio Emergenza/Urgenza 118.

Le domande poste ai dipendenti della ASL Foggia sono identiche a quelle utilizzate per l'edizione 2018. Per semplificare la distribuzione dei questionari si è preferito creare dei questionari ad hoc per ogni figura rispondente. Quindi per ogni rispondente del nostro frame abbiamo realizzato un questionario "su misura" in base alle attività che doveva valutare.

I questionari adottati per la rilevazione della customer satisfaction sono composti da quattro sezioni:

- Sezione 1 che contiene una breve descrizione della finalità dell'indagine, le informazioni relative alla struttura di riferimento ed il ruolo che il rispondente ricopre all'interno della stessa.

- Sezione 2, contenente le informazioni relative alle diverse attività da valutare. Le attività da valutare variano in base ai servizi analizzati e sono state indivi-

duate sintetizzando la descrizione molto puntuale presente nel capitolato stipulato con la ASL di Foggia.

Ogni questionario considera una o più attività a seconda dei servizi di cui è responsabile il rispondente.

Per ciascuna attività analizzata gli intervistati hanno indicato quattro distinti punteggi da 1 a 10 a seconda degli aspetti da valutare:

- o professionalità
- o tempo di intervento
- o qualità dell'intervento
- o comportamento capacità relazionali
- o soddisfazione complessiva
- Sezione 3, dedicata all'approfondimento dei diversi aspetti specifici dei servizi considerati, al livello di conoscenza delle procedure del personale ed alle attività di formazione realizzate da Sanitaservice ASL FG.
- Sezione 4, dedicata allo studio della rete di relazioni interne a Sanitaservice ASL FG al fine di capire quali siano i servizi che hanno più frequenti interazioni. Il questionario si conclude con una domanda relativa alla valutazione complessiva dell'organizzazione di Sanitaservice ASL FG e con uno spazio a risposta aperta destinato a eventuali commenti e suggerimenti.

2.2 La popolazione oggetto dell'indagine

La precedente edizione dell'indagine, effettuata nel 2018, aveva coinvolto ben 129 intervistati che a seconda del numero dei servizi loro affidati avevano compilato 260 questionari.

Nell'edizione 2019 i questionari raccolti sono stati 317 ma gli intervistati sono stati solo 103. I motivi di tali scostamenti sono riconducibili a:

- la quasi assenza nella rilevazione del Presidio Ospedaliero "F. Lastaria" di Lucera, in quanto passato durante la seconda metà del 2019 sotto la giurisdizione degli Operali Riuniti di Foggia e quindi uscendo dalla competenza dell'ASL FG;
- l'assenza di diversi Direttori Sanitari di Dipartimenti di Aree e Servizi dell'ASL FG, in quanto non più fruitori dei servizi erogati dalla Sanitaservice ASL FG;
- la riorganizzazione di molti reparti all'interno dei vari presidi ospedalieri che ha comportato una consistente riduzione della unità operative partecipanti alla rilevazione;

- l'aumento del numero dei questionari relativi all'Attività di Pulizie e di Manutenzione rispetto alle precedenti edizioni.

Possiamo quindi affermare che, nel periodo 5 marzo 2020 - 23 giugno 2020, i questionari raccolti dai responsabili della Sanitaservice ASL FGASL FG sono stati pari al 94,5% di quelli previsti. Tenuto conto del periodo di grande emergenza e di superlavoro dovuto alla emergenza Covid-19 in cui si è svolta la rilevazione, l'ampio tasso di copertura raggiunto è ampiamente soddisfacente.

Tabella 1. *Numero questionari raccolti relativi ai diversi servizi erogati da Sanitaservice ASL FG nel 2019. Rilevazione 2017/2018/2019.*

SERVIZI EROGATI DA SANITASERVICE ASL FG	QUESTIONARI RACCOLTI		
	2017	2018	2019
<i>Servizio di Ausiliariato</i>	14	101	94
<i>Servizio di Pulizie e Sanificazione ambientale</i>	15	58	96
<i>Servizio di Manutenzione degli immobili</i>	13	77	96
<i>Servizio di Trasporto beni economici</i>	12	13	14
<i>Servizio di Ausiliariato e Infermieristico UU.DD.TT.</i>	6	6	11
<i>Servizio Emergenza/Urgenza 118</i>	3	5	6
TOTALE	63	260	317

3. I risultati della Customer Satisfaction

3.1 Servizio di Ausiliariato

I rispondenti hanno assegnato un punteggio da 1 a 10 alla propria percezione delle seguenti quattro dimensioni da valutare: Professionalità, Comportamento, Capacità relazionali e Soddisfazione complessiva.

La soddisfazione media complessiva è stata pari a 9,0, ben 0,5 punti in più rispetto all'edizione precedente dove il valore medio della soddisfazione è stato pari a 8,5 e 1,1 punti in più rispetto alla rilevazione 2017. Tali valori indicano chiaramente che le politiche attuate dalla Sanitaservice ASL FG in questi anni stanno producendo i risultati attesi.

Una particolarità da segnalare riguarda l'attività di trasporto di medicine, referti, materiale biologico, ecc.... Come possiamo osservare il servizio è migliorato negli anni e si è sempre mantenuto su livelli particolarmente elevati sin dalla prima rilevazione effettuata.

Anche i valori delle altre 3 dimensioni valutate sono molto più alti rispetto al 2018 e al 2017. Ciò dimostra che il miglioramento delle performances è stato generalizzato ed è attribuibile ad una crescita del livello di professionalità e della capacità relazionali ma anche ad un sensibile miglioramento dei comportamenti posti in essere dal personale valutato.

Tabella 2. *Punteggi in scala 1 – 10 attribuiti dai rispondenti ai diversi servizi erogati da Sanitaservice ASL FG relativi all’attività di ausiliario. Rilevazione 2019/2018/2017.*

FREQ. ASSOL.	ATTIVITÀ DA VALUTARE AUSILIARIATO anno 2019	ASPETTI DA VALUTARE			
		PROF.	COMP.	REL.	SOD.
94	1. Accompagnamento e trasporto dei degenti con mezzi adeguati	9,1	9,1	9,1	9,1
94	2. Pulizia degli ambienti e operazioni elementari e di supporto necessarie al funzionamento del reparto, quali lo spostamento dei ricoverati	9,1	9,1	9,1	9,1
94	3. Trasporto di medicine, referti, materiale biologico, sanitario ed economale, vitto, attrezzature, vestiario, biancheria, etc.	9,1	9,1	9,0	9,1
94	4. Areare, spazzare, lavare e spolverare tutti gli ambienti dell’unità operativa alla quale è addetto	9,0	9,0	9,0	9,0
94	5. Partecipazione con l’équipe di lavoro, limitatamente ai propri compiti	8,9	9,0	8,9	8,9
94	6. Controllo degli accessi delle persone	8,8	8,8	8,8	8,8
	MEDIA	9,0	9,0	9,0	9,0
FREQ. ASSOL.	ATTIVITÀ DA VALUTARE AUSILIARIATO anno 2018	PROF.	COMP.	REL.	SOD.
101	1. Accompagnamento e trasporto dei degenti con mezzi adeguati	8,3	8,4	8,4	8,4
101	2. Pulizia degli ambienti e operazioni elementari e di supporto necessarie al funzionamento del reparto, quali lo spostamento dei ricoverati	8,6	8,6	8,5	8,6
101	3. Trasporto di medicine, referti, materiale biologico, sanitario ed economale, vitto, attrezzature, vestiario, biancheria, etc.	8,6	8,7	8,6	8,8
101	4. Areare, spazzare, lavare e spolverare tutti gli ambienti dell’unità operativa alla quale è addetto	8,6	8,6	8,6	8,5
101	5. Partecipazione con l’équipe di lavoro, limitatamente ai propri compiti	8,5	8,6	8,5	8,6
101	6. Controllo degli accessi delle persone	8,3	8,3	8,3	8,3
	MEDIA	8,5	8,5	8,5	8,5

Segue Tabella 2. *Punteggi in scala 1 – 10 attribuiti dai rispondenti ai diversi servizi erogati da Sanitaservice ASL FG relativi all'attività di ausiliario. Rilevazione 2019/2018/2017.*

FREQ. ASSOL.	ATTIVITÀ DA VALUTARE AUSILIARIATO anno 2017	ASPETTI DA VALUTARE			
		PROF.	COMP.	REL.	SOD.
15	1. Accompagnamento e trasporto dei degenti con mezzi adeguati	7,9	7,9	7,7	7,9
15	2. Pulizia degli ambienti e operazioni elementari e di supporto necessarie al funzionamento del reparto, quali lo spostamento dei ricoverati	7,9	7,8	7,7	7,7
15	3. Trasporto di medicine, referti, materiale biologico, sanitario ed economale, vitto, attrezzature, vestiario, biancheria, etc.	8,4	8,5	8,4	8,3
15	4. Areare, spazzare, lavare e spolverare tutti gli ambienti dell'unità operativa alla quale è addetto	7,9	8,1	7,8	7,9
15	5. Partecipazione con l'équipe di lavoro, limitatamente ai propri compiti	7,9	7,9	7,8	7,9
15	6. Controllo degli accessi delle persone	7,7	7,8	7,8	7,7
	MEDIA	8,0	8,0	7,9	7,9

Come si vede dalla figura 1, in cui sono riportate le medie de punteggi attribuiti alle diverse attività valutate, l'incremento dal 2017 al 2019 è stato di circa mezzo punto all'anno per tutti gli aspetti considerati.

Figura 1. *Confronto fra le medie dei diversi aspetti da valutare relativi alle attività di ausiliario per gli anni 2019/2018/2017.*



3.2 Servizio di Pulizia e Sanificazione ambientale

In questa sezione i rispondenti hanno dovuto valutare quattro dimensioni, ovvero: Professionalità, Tempestività d'intervento, Qualità dell'intervento e Soddisfazione complessiva; sempre utilizzando una scala di valutazione che va da 1 a 10.

La soddisfazione complessiva media pari a 9,0 punti è aumentata rispetto agli 8,1 punti registrati l'anno precedente ed ai 7,8 punti del 2017.

Tabella 3. *Punteggi attribuiti ai diversi servizi erogati da SS all'attività di pulizia. Rilevazione 2019/2018/2017.*

FREQ. ASSOL.	ATTIVITÀ DA VALUTARE: PULIZIE anno 2019	ASPETTI DA VALUTARE			
		PROF.	COMP.	REL.	SOD.
96	1. Attività ordinarie di pulizia e sanificazione e di- sinfezione	9,2	9,0	9,1	9,1
96	2. Decontaminazione e rimozione di eventuale ma- teriale organico da tutte le superfici	9,1	9,0	9,0	9,0
96	3. Raccolta e trasporto di tutte le categorie di ri- fiu- ti, sostituzione del sacchetto, deterzione e di- sinfe- zione dei contenitori	9,1	9,0	9,1	9,1
96	4. Pulizia e disinfezione degli arredi e delle attrez- zature mobili	9,0	9,0	9,0	9,0
96	5. Spolveratura ad umido, deterzione e successiva disinfezione, di tutte le superfici orizzontali e ver- ticali (altezza 180 cm)	9,0	8,9	8,9	9,0
96	6. Deterzione e disinfezione dei lavabi, accessori e arredi sanitari	9,0	9,0	9,0	9,0
96	7. Pulizia a fondo e disinfezione degli arredi mo- bi- li (carrelli, tavoli con ruote, ecc.)	8,9	8,8	8,8	8,9
96	8. Disincrostazione sanitari, rubinetterie e zone li- mitrofe	8,8	8,8	8,8	8,8
96	9. Deragnatura	9,0	8,8	8,9	8,9
96	10. Raccolta e trasporto dei rifiuti speciali fino al punto di deposito temporaneo	9,0	8,9	9,0	9,0
96	11. Palestre: deterzione materassini, parallele, sca- le, ecc.	8,9	8,9	8,9	8,9
	MEDIA	9,0	8,9	9,0	9,0

Segue Tabella 3. *Punteggi attribuiti ai diversi servizi erogati da SS all'attività di pulizia. Rilevazione 2019/2018/2017.*

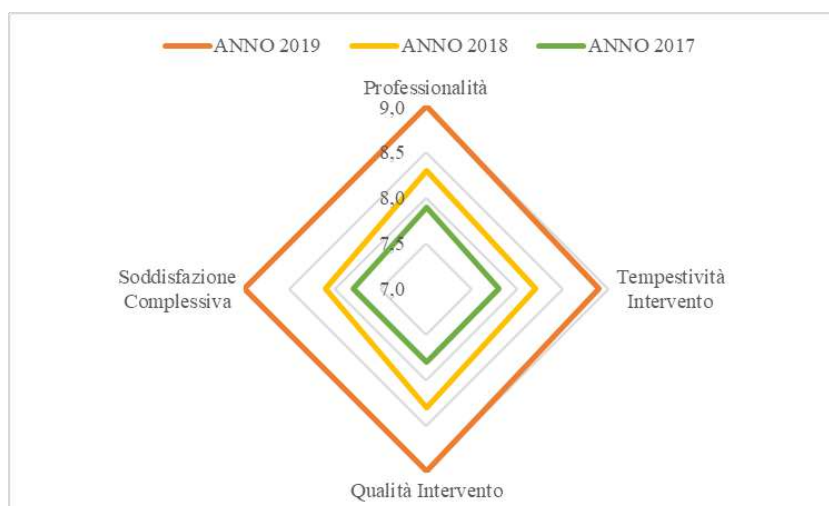
FREQ. AS- SOL.	ATTIVITÀ DA VALUTARE: PULIZIE anno 2018	ASPETTI DA VALUTARE			
		PROF.	COMP.	REL.	SOD.
58	1. Attività ordinarie di pulizia e sanificazione e disinfezione	8,3	8,4	8,4	8,4
58	2. Decontaminazione e rimozione di eventuale materiale organico da tutte le superfici	8,3	8,3	8,3	8,2
58	3. Raccolta e trasporto di tutte le categorie di rifiuti, sostituzione del sacchetto, deterzione e disinfezione dei contenitori	8,4	8,3	8,4	8,3
58	4. Pulizia e disinfezione degli arredi e delle attrezzature mobili	8,2	8,2	8,2	8,1
58	5. Spolveratura ad umido, deterzione e successiva disinfezione, di tutte le superfici orizzontali e verticali (altezza 180 cm)	8,2	8,1	8,2	7,9
58	6. Deterzione e disinfezione dei lavabi, accessori e arredi sanitari	8,3	8,3	8,3	8,1
58	7. Pulizia a fondo e disinfezione degli arredi mobili (carrelli, tavoli con ruote, ecc.)	8,1	8,0	8,1	7,9
58	8. Disincrostazione sanitari, rubinetterie e zone limitrofe	8,1	8,1	8,1	8,0
58	9. Deragnatura	7,9	8,0	8,0	8,0
58	10. Raccolta e trasporto dei rifiuti speciali fino al punto di deposito temporaneo	8,4	8,4	8,5	8,3
58	11. Palestre: deterzione materassini, parallele, scale, ecc.	8,3	8,4	8,4	8,3
	MEDIA	8,3	8,2	8,3	8,1
FREQ. AS- SOL.	ATTIVITÀ DA VALUTARE: PULIZIE anno 2017	ASPETTI DA VALUTARE			
		PROF.	COMP.	REL.	SOD.
18	1. Attività ordinarie di pulizia e sanificazione e disinfezione	8,1	8,0	8,1	8,1
18	2. Decontaminazione e rimozione di eventuale materiale organico da tutte le superfici	8,1	7,9	8,0	8,0
18	3. Raccolta e trasporto di tutte le categorie di rifiuti Sostituzione del sacchetto, deterzione e disinfezione dei contenitori	8,1	7,9	8,1	8,0
18	4. Pulizia e disinfezione degli arredi e delle attrezzature mobili	7,6	7,4	7,5	7,5

Segue Tabella 3. *Punteggi attribuiti ai diversi servizi erogati da SS all'attività di pulizia. Rilevazione 2019/2018/2017.*

FREQ. ASSOL.	ATTIVITÀ DA VALUTARE: PULIZIE anno 2017	ASPETTI DA VALUTARE			
		PROF.	COMP.	REL.	SOD.
18	5. Spolveratura ad umido, deterzione e successiva disinfezione, di tutte le superfici orizzontali e verticali (altezza 180 cm)	7,4	7,3	7,4	7,4
18	6. Deterzione e disinfezione dei lavabi, accessori e arredi sanitari	7,8	7,8	7,7	7,8
18	7. Pulizia a fondo e disinfezione degli arredi mobili (carrelli, tavoli con ruote, ecc.)	7,8	7,6	7,7	7,7
18	8. Disincrostazione sanitari, rubinetterie e zone limitrofe	7,6	7,5	7,5	7,5
18	9. Deragnatura	7,4	7,4	7,4	7,4
18	10. Raccolta e trasporto dei rifiuti speciali fino al punto di deposito temporaneo	8,3	8,3	8,3	8,3
18	11. Palestre: deterzione materassini, ecc. (esclusi i giocattoli)	8,1	8,1	7,9	8,0
	MEDIA	7,9	7,8	7,8	7,8

Come si evince dalla tabella 3 tutte le attività valutate hanno registrato dei miglioramenti. È importante sottolineare come il valore più basso del 2019 pari a 8,8 punti è sensibilmente più alto del voto massimo del 2018 risultato pari a 8,4 e del 2017 pari a 8,3.

Figura 2. *Confronto fra i diversi aspetti da valutare relativi alle attività di pulizia per gli anni 2019/2018/2017.*



Analizzando le medie complessive delle dimensioni attraverso la figura 2 si ha una conferma immediata della crescita registrata rispetto agli anni precedenti.

3.3 Servizio di Manutenzione degli immobili

Anche per questa tipologia di servizio i rispondenti hanno dovuto valutare 4 distinte dimensioni: Professionalità, Tempestività d'intervento, Qualità dell'intervento e la Soddisfazione complessiva; utilizzando una scala di punteggi discreti 1-10.

La soddisfazione media risulta pari a 8,5 punti, leggermente maggiore rispetto al valore medio di 8,2 registrato nel 2018 e al valore 7,3 del 2017.

In coerenza con quanto riscontrato nelle edizioni precedenti, il servizio che ha avuto la valutazione più bassa è quello relativo alla manutenzione del verde. Evidentemente le strutture sanitarie prestano più attenzione alla cura dei loro assistiti piuttosto che agli spazi verdi annessi alle strutture valutate.

Tabella 4. *Punteggi in scala 1 – 10 attribuiti dai rispondenti ai diversi servizi erogati da Sanitaservice ASL FG relativi all'attività di manutenzione. Rilevazione 2019/2018/2017.*

FREQ. ASSOL.	ATTIVITÀ DA VALUTARE MANUTENZIONE anno 2019	ASPETTI DA VALUTARE			
		PROF	COMP.	REL.	SOD.
96	1. Gli interventi di manutenzione effettuati da parte di Sanitaservice nei locali interni ed esterni (a titolo esemplificativo e non esaustivo: interventi su infissi, opere murarie, cancellate, impianti di scarico delle acque, etc.)	8,8	8,5	8,7	8,7
96	2. Mantenere nelle condizioni estetiche migliori l'arredo verde esterno, intesa qualunque coltura arborea e floreale	8,4	8,2	8,2	8,3
96	3. Gestione e smistamento delle richieste di manutenzione ordinaria sugli impianti e gli immobili, pervenuti sui protocolli informatici dei vari presidi ospedalieri e dai D.S.S. dell'ASL di Foggia	8,5	8,4	8,5	8,4
96	4. Supporto alla pianificazione e programmazione degli interventi di manutenzione da eseguire presso le varie strutture sanitarie dell'ASL	8,7	8,5	8,7	8,6
	MEDIA	8,6	8,4	8,5	8,5

Segue Tabella 4. *Punteggi in scala 1 – 10 attribuiti dai rispondenti ai diversi servizi erogati da Sanitaservice ASL FG relativi all'attività di manutenzione. Rilevazione 2019/2018/2017.*

FREQ. ASSOL.	ATTIVITÀ DA VALUTARE MANUTENZIONE anno 2018	ASPETTI DA VALUTARE			
		PROF	COMP.	REL.	SOD.
77	1. Gli interventi di manutenzione effettuati da parte di Sanitaservice nei locali interni ed esterni (a titolo esemplificativo e non esaustivo: interventi su infissi, opere murarie, cancellate, impianti di scarico delle acque, etc.)	8,5	8,4	8,4	8,3
77	2. Mantenere nelle condizioni estetiche migliori l'arredo verde esterno, intesa qualunque coltura arborea e floreale	8,0	8,0	8,0	7,9
77	3. Gestione e smistamento delle richieste di manutenzione ordinaria sugli impianti e gli immobili, pervenuti sui protocolli informatici dei vari presidi ospedalieri e dai D.S.S. dell'ASL di Foggia	8,3	8,2	8,3	8,3
77	4. Supporto alla pianificazione e programmazione degli interventi di manutenzione da eseguire presso le varie strutture sanitarie dell'ASL	8,1	8,1	8,1	8,2
	MEDIA	8,3	8,2	8,2	8,2
FREQ. ASSOL.	ATTIVITÀ DA VALUTARE MANUTENZIONE anno 2017	ASPETTI DA VALUTARE			
		PROF	COMP.	REL.	SOD.
12	1. Gli interventi di manutenzione effettuati da parte di Sanitaservice nei locali interni ed esterni	7,8	8,2	7,9	7,9
12	2. Mantenere nelle condizioni estetiche migliori l'arredo verde esterno, intesa qualunque coltura arborea e floreale	7,3	7,1	6,9	6,9
12	3. Gestione e smistamento delle richieste di manutenzione ordinaria sugli impianti e gli immobili, pervenuti sui protocolli informatici dei vari presidi ospedalieri e dai D.S.S. dell'ASL di Foggia	7,4	7,4	7,4	7,4
12	4. Supporto alla pianificazione e programmazione degli interventi di manutenzione da eseguire presso le varie strutture sanitarie dell'ASL	7,3	7,2	7,1	7,2
	MEDIA	7,4	7,5	7,3	7,3

La figura 3 ci fornisce una idea immediata dei miglioramenti medi registrati nei 3 anni considerati.

Figura 3. Confronto fra i diversi aspetti da valutare relativi alle attività di manutenzione per gli anni 2019/2018/2017.



Gli intervistati sulle attività di manutenzione hanno risposto anche a delle domande di approfondimento per capire se conoscono quali siano le attività di manutenzione di competenza di Sanitaservice ASL FG e quelle affidate a ditte esterne. Come si vede dalla figura 4, l'86,6% degli intervistati conosce bene i servizi di competenza dei manutentori di Sanitaservice ASL FG. I rispondenti sono, inoltre, convinti che gli addetti ai servizi di manutenzione si adoperino al meglio per minimizzare i disagi dei clienti durante i loro interventi, come si vede dalla figura 5.

Figura 4. Risposte fornite dagli intervistati relative all'esistenza di una chiara distinzione tra le attività di competenza della Sanitaservice ASL FG e le attività di manutenzione affidate a ditte esterne.

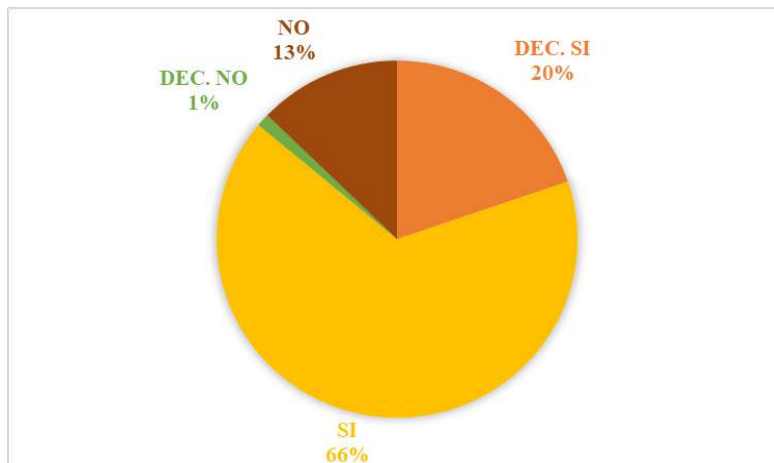
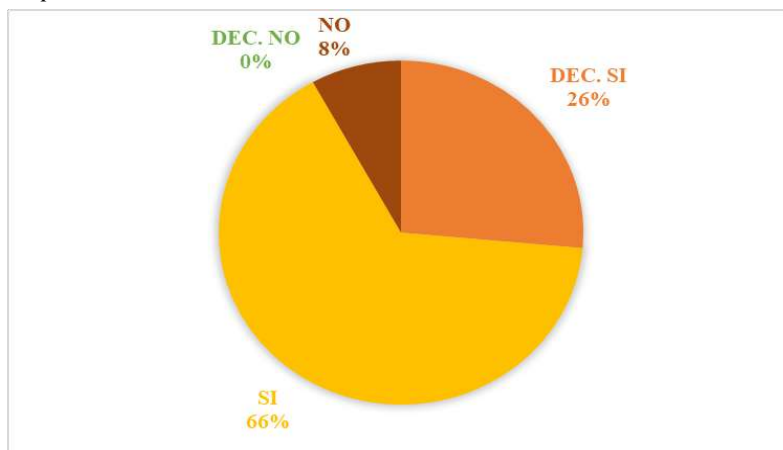


Figura 5. Risposte fornite dagli intervistati relative alla capacità di Sanitaservice ASL FG di effettuare i lavori di manutenzione minimizzando i disagi agli utenti e agli operatori delle strutture.



3.4 Servizio di Trasporto beni economici

Nell'analizzare questa specifica attività, non notiamo grandissime novità; rileviamo una soddisfazione media complessiva pari 8,8 rispetto all'8,6 del 2018 e all'8,3 del 2017.

Possiamo osservare come l'attività di Movimentazione dei flussi documentali, di materiale biologico e dei referti tra le strutture dell'ASL di Foggia, sia rimasta pressoché invariata nei giudizi, a differenza delle altre due attività che invece hanno ottenuto valutazioni leggermente maggiori rispetto agli anni precedenti.

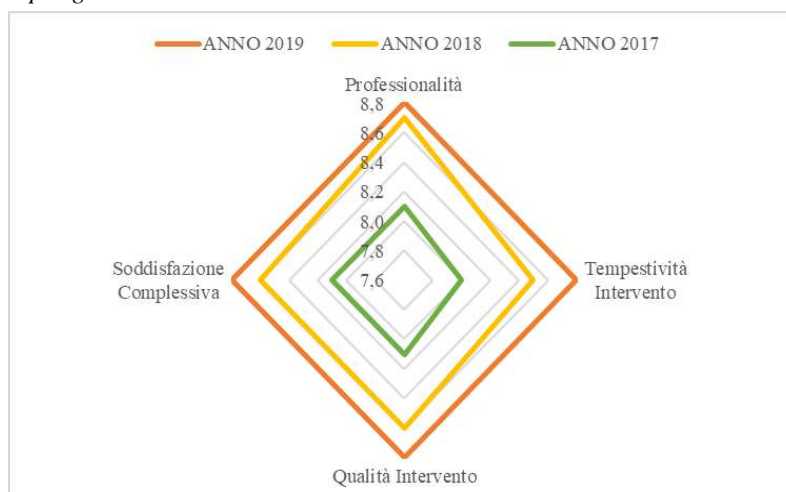
Tabella 5. Punteggi in scala 1 – 10 attribuiti dai rispondenti ai diversi servizi erogati da Sanitaservice ASL FG relativi all'attività di trasporto beni economici. Rilevazione 2019/2018/2017.

FREQ. ASSOL.	ATTIVITÀ DA VALUTARE: TRASPORTO BENI ECONOMICI anno 2019	ASPETTI DA VALUTARE			
		PROF.	COMP.	REL.	SOD.
14	1. Movimentazione dei flussi documentali, di materiale biologico e dei referti tra le strutture dell'ASL di Foggia	8,6	8,5	8,6	8,6
14	2. Prelievo, trasporto e distribuzione farmaci, materiale sanitario e cancelleria	8,9	8,9	9,0	8,9
14	3. Prelievo, trasporto e consegna di ricette farmaceutiche e specialistiche	9,0	8,9	8,9	8,9
	MEDIA	8,8	8,8	8,8	8,8

Segue Tabella 5. *Punteggi in scala 1 – 10 attribuiti dai rispondenti ai diversi servizi erogati da Sanitaservice ASL FG relativi all'attività di trasporto beni economici. Rilevazione 2019/2018/2017.*

FREQ. ASSOL.	ATTIVITÀ DA VALUTARE: TRASPORTO BENI ECONOMICI anno 2018	ASPETTI DA VALUTARE			
		PROF	COMP.	REL.	SOD.
13	1. Movimentazione dei flussi documentali, di materiale biologico e dei referti tra le strutture dell'ASL di Foggia	8,8	8,5	8,6	8,6
13	2. Prelievo, trasporto e distribuzione farmaci, materiale sanitario e cancelleria	8,6	8,5	8,5	8,5
13	3. Prelievo, trasporto e consegna di ricette farmaceutiche e specialistiche	8,6	8,5	8,6	8,7
	MEDIA	8,7	8,5	8,6	8,6
FREQ. ASSOL.	ATTIVITÀ DA VALUTARE: TRASPORTO BENI ECONOMICI anno 2017	ASPETTI DA VALUTARE			
		PROF	COMP.	REL.	SOD.
12	1. Movimentazione dei flussi documentali, di materiale biologico e dei referti tra le strutture dell'ASL di Foggia	8,2	8,1	8,2	8,2
12	2. Prelievo, trasporto e distribuzione farmaci, materiale sanitario e cancelleria	8,1	8,1	8,2	8,1
12	3. Prelievo, trasporto e consegna di ricette farmaceutiche e specialistiche	7,9	7,9	7,9	7,9
	MEDIA	8,1	8,0	8,1	8,1

Figura 6. *Confronto fra i diversi aspetti da valutare relativi alle attività di trasporto beni economici per gli anni 2019/2018/2017.*



Nel questionario erano presenti anche due domande di approfondimento:

- Ritiene che il personale di Sanitaservice ASL FG rispetti le procedure atte a garantire la riservatezza su documenti, informazioni e altro materiale?
- Ritiene che il personale di Sanitaservice ASL FG rispetti il divieto di fumare negli spazi interdetti?

Figura 7. Risposte fornite dagli intervistati relative al rispetto delle procedure di riservatezza da parte degli addetti della Sanitaservice ASL FG.

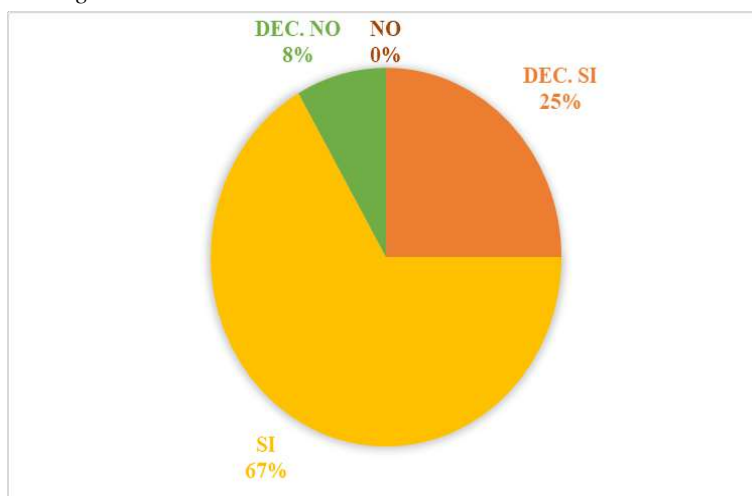
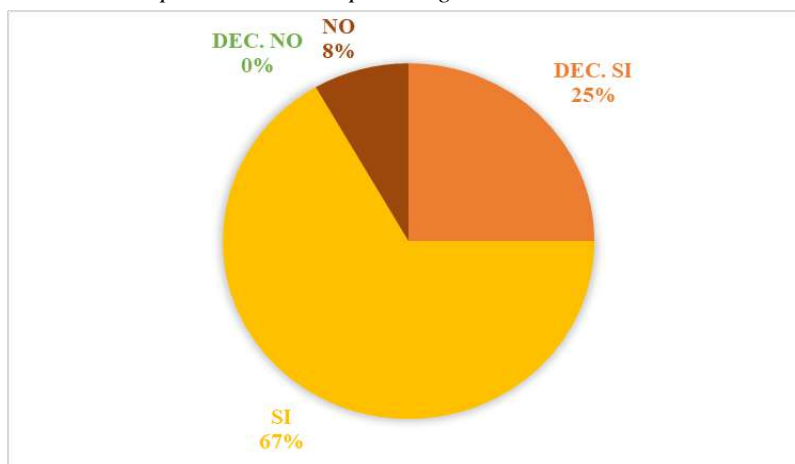


Figura 8. Risposte fornite dagli intervistati relative al rispetto di divieto di fumare negli spazi interdetti da parte degli addetti Sanitaservice ASL FG.



Come si vede dalle figure 7 e 8, le risposte fornite dagli intervistati indicano che i comportamenti degli addetti al trasporto dei beni economici sono virtuosi.

3.5 Servizio di Ausiliariato dell'UU.DD.TT.

Le Unità di Degenza Territoriale (UU.DD.TT.) costituiscono l'alternativa all'assistenza domiciliare integrata laddove non ci sia una famiglia in grado di supportare il personale sanitario.

Tabella 6. *Punteggi in scala 1 – 10 attribuiti dai rispondenti ai diversi servizi erogati da Sanitaservice ASL FG relativi all'attività di ausiliariato UU.DD.TT.. Rilevazione 2019/2018/2017.*

FREQ. ASSOL.	ATTIVITÀ DA VALUTARE AUSILIARIATO UU.DD.TT. anno 2019	ASPETTI DA VALUTARE			
		PROF	COMP.	REL.	SOD.
11	1. Accompagnamento e trasporto dei degenti con mezzi adeguati	7,3	7,1	7,1	7,3
11	2. Pulizia degli ambienti e operazioni elementari e di supporto necessarie al funzionamento del reparto quali lo spostamento dei ricoverati	8,3	8,4	8,4	8,5
11	3. Pulizia dell'ambiente dopo il pasto e aiuto nella distribuzione e somministrazione del pasto	8,4	8,4	8,4	8,5
11	4. Collaborazione con l'infermiere professionale per atti di accudimento semplici al malato	8,3	8,3	8,3	8,4
11	5. Areare, spazzare, lavare e spolverare tutti gli ambienti dell'unità operativa alla quale è addetto	8,1	8,0	7,9	8,3
11	6. Pulizia e manutenzione di utensili, apparecchi, presidi usati dal paziente e dal personale medico ed infermieristico per l'assistenza del malato	7,9	8,0	8,0	8,0
11	7. Comunicazione all'infermiere professionale di quanto sopravviene durante il suo lavoro in quanto ritenuto in quanto ritenuto incidente sull'assistito e sull'ambiente	7,9	7,9	7,9	7,9
11	8. Partecipazione con l'equipe di lavoro, limitatamente ai propri compiti	7,8	7,9	7,8	7,9
11	9. Rifacimento del letto non occupato (comodino, letto e apparecchiature)	8,5	8,5	8,5	8,5
	MEDIA	8,1	8,1	8,0	8,1
FREQ. ASSOL.	ATTIVITÀ DA VALUTARE AUSILIARIATO UU.DD.TT. anno 2018	ASPETTI DA VALUTARE			
		PROF	COMP.	REL.	SOD.
6	1. Accompagnamento e trasporto dei degenti con mezzi adeguati	7,0	7,0	7,2	6,8
6	2. Pulizia degli ambienti e operazioni elementari e di supporto necessarie al funzionamento del reparto quali lo spostamento dei ricoverati	6,8	7,2	7,2	7,0

Segue Tabella 6. Punteggi in scala 1 – 10 attribuiti dai rispondenti ai diversi servizi erogati da Sanitaservice ASL FG relativi all'attività di ausiliariato UU.DD.TT.. Rilevazione 2019/2018/2017.

FREQ. ASSOL.	ATTIVITÀ DA VALUTARE AUSILIARIATO UU.DD.TT. anno 2018	ASPETTI DA VALUTARE			
		PROF	COMP.	REL.	SOD.
6	3. Pulizia dell'ambiente dopo il pasto e aiuto nella distribuzione e somministrazione del pasto	7,3	7,2	7,3	7,2
6	4. Collaborazione con l'infermiere professionale per atti di accudimento semplici al malato	7,0	6,8	6,8	6,8
6	5. Areare, spazzare, lavare e spolverare tutti gli ambienti dell'unità operativa alla quale è addetto	7,0	7,0	7,0	7,0
6	6. Pulizia e manutenzione di utensili, apparecchi, presidi usati dal paziente e dal personale medico ed infermieristico per l'assistenza del malato	6,8	7,0	7,0	7,0
6	7. Comunicazione all'infermiere professionale di quanto sopravviene durante il suo lavoro in quanto ritenuto incidente sull'assistito e sull'ambiente	7,0	7,0	7,0	7,0
6	8. Partecipazione con l'equipe di lavoro, limitatamente ai propri compiti	7,0	7,0	7,2	6,8
6	9. Rifacimento del letto non occupato (comodino, letto e apparecchiature)	7,3	7,3	7,3	7,3
	MEDIA	7,0	7,1	7,1	7,0
FREQ. ASSOL.	ATTIVITÀ DA VALUTARE AUSILIARIATO UU.DD.TT. anno 2017	PROF	COMP.	REL.	SOD.
6	3. Pulizia dell'ambiente dopo il pasto e aiuto nella distribuzione e somministrazione del pasto	6,8	7,0	7,0	6,8
6	4. Collaborazione con l'infermiere professionale per atti di accudimento semplici al malato	7,0	7,0	6,8	7,0
6	6. Pulizia a manutenzione di utensili, apparecchi, presidi usati dal paziente e dal personale medico ed infermieristico per l'assistenza al malato	7,0	7,0	7,0	7,0
6	7. Comunicazione all'infermiere professionale di quanto sopravviene durante il suo lavoro in quanto ritenuto incidente sull'assistito e sull'ambiente	7,0	7,0	6,8	7,0
6	9. Rifacimento del letto non occupato (comodino, letto e apparecchiature)	7,2	7,2	7,0	7,2
	MEDIA	7,0	7,0	6,9	7,0

In generale tutte le altre attività sono decisamente migliorate. La soddisfazione complessiva riferita al 2019 è pari ad 8,1 ben superiore del 7,0 registrato nel 2018, tenendo conto che proprio nel 2018 non c'era stato un grande miglioramento rispetto al 2017; ciò implica come tale servizio di ausiliario UU.DD.TT. sia stato oggetto di lavoro da parte della Sanitaservice ASL FG permettendo di arrivare ad una soddisfazione largamente migliore rispetto a quella rilevata in passato.

Figura 9. *Confronto fra i diversi aspetti da valutare relativi alle attività di ausiliario UU.DD.TT. per gli anni 2019/2018/2017.*



3.6 Servizio Infermieristico dell'UU.DD.TT.

Il servizio infermieristico dell'UU.DD.TT. è sicuramente quello che in termini differenziali ha avuto la maggior crescita passando da una soddisfazione complessiva pari a 7 nel 2018 ad una pari a 8,5 nel 2019 recuperando in parte anche rispetto al voto 7,3 del 2017.

Come negli anni precedenti l'attività: Compilazione della cartella clinica (per quanto attiene la parte infermieristica) è quella con la valutazione più bassa rispetto alle altre ma, nel 2019, ci sono stati degli incrementi importanti rispetto al 2018 visto che il punteggio medio assegnato è passato a 8,0 rispetto al 6,6 dell'edizione scorsa.

La figura 10 conferma che la Sanitaservice ASL FG sta migliorando i propri servizi rispetto a tutte e 4 le dimensioni valutate.

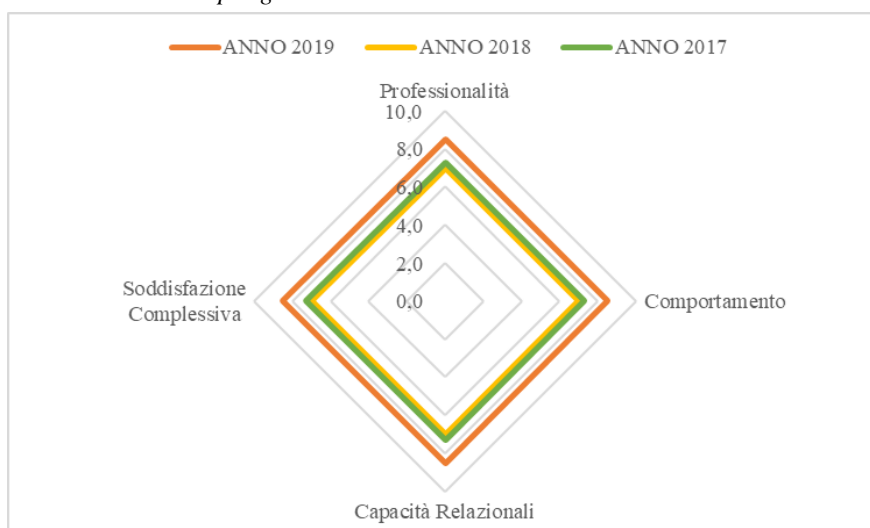
Tabella 7. *Punteggi in scala 1 – 10 attribuiti dai rispondenti ai diversi servizi erogati da Sanitaservice ASL FG relativi all'attività infermiere UU.DD.TT.. Rilevazione 2019/2018/2017.*

FREQ. ASSOL.	ATTIVITÀ DA VALUTARE: INFERMIERE UU.DD.TT. anno 2019	ASPETTI DA VALUTARE			
		PROF.	COMP.	REL.	SOD.
11	1. Accoglienza dei ricoverati e dei loro familiari informandoli sulle caratteristiche della struttura e dell'organizzazione assistenziale	8,4	8,1	8,1	8,1
11	2. Compilazione della cartella clinica per quanto attiene alla parte infermieristica	7,9	7,9	8,0	7,9
11	3. Vigilanza sullo stato del paziente	8,6	8,5	8,5	8,6
11	4. Esecuzione del programma assistenziale programmato dal MMG collaborando attivamente nel raggiungimento degli obiettivi di salute	8,6	8,6	8,6	8,6
11	5. Controllo dei parametri vitali	8,6	8,9	8,6	8,8
11	6. Controllo dell'igiene e profilassi anti-decubito	8,3	8,3	8,3	8,4
11	7. Rifacimento letti (con paziente allettato) e igiene dell'unità di vita del paziente	8,3	8,3	8,3	8,3
11	8. Prelievo sangue	8,9	8,9	8,9	8,9
11	9. Rilevamento temperatura corporea	8,8	8,9	8,9	8,9
	MEDIA	8,5	8,5	8,5	8,5
FREQ. ASSOL.	ATTIVITÀ DA VALUTARE: INFERMIERE UU.DD.TT. anno 2018	ASPETTI DA VALUTARE			
		PROF.	COMP.	REL.	SOD.
6	1. Accoglienza dei ricoverati e dei loro familiari informandoli sulle caratteristiche della struttura e dell'organizzazione assistenziale	6,8	6,8	6,8	6,8
6	2. Compilazione della cartella clinica per quanto attiene alla parte infermieristica	6,7	6,7	6,5	6,5
6	3. Vigilanza sullo stato del paziente	7,2	7,2	7,2	7,2
6	4. Esecuzione del programma assistenziale programmato dal MMG collaborando attivamente nel raggiungimento degli obiettivi di salute	6,8	6,8	6,7	6,8
6	5. Controllo dei parametri vitali	7,2	7,2	7,2	7,2
6	6. Controllo dell'igiene e profilassi anti-decubito	6,8	7,0	6,8	6,8
6	7. Rifacimento letti (con paziente allettato) e igiene dell'unità di vita del paziente	7,2	7,2	7,2	7,0
6	8. Prelievo sangue	7,2	7,2	7,2	7,2
6	9. Rilevamento temperatura corporea	7,2	7,2	7,2	7,2
	MEDIA	7,0	7,0	7,0	7,0

Segue Tabella 7. Punteggi in scala 1 – 10 attribuiti dai rispondenti ai diversi servizi erogati da Sanitaservice ASL FG relativi all'attività infermiere UU.DD.TT.. Rilevazione 2019/2018/2017.

FREQ. ASSOL.	ATTIVITÀ DA VALUTARE: INFERMIERE UU.DD.TT. anno 2017	ASPETTI DA VALUTARE			
		PROF	COMP	REL.	SOD.
6	1. Accoglienza dei ricoverati e dei loro familiari informandoli sulle caratteristiche della struttura e dell'organizzazione assistenziale	7,0	6,8	6,8	6,8
6	2. Compilazione della cartella clinica per quanto attiene alla parte infermieristica	6,7	6,8	6,8	6,8
6	3. Vigilanza sullo stato del paziente	7,2	7,2	7,2	7,2
6	4. Esecuzione del programma assistenziale programmato dal MMG collaborando attivamente nel raggiungimento degli obiettivi di salute	7,0	7,0	6,7	6,8
6	5. Controllo dei parametri vitali	7,8	7,8	7,7	7,8
6	6. Controllo dell'igiene e profilassi anti-decubito	7,3	7,3	7,3	7,2
6	7. Rifacimento letti (con paziente allettato) e igiene dell'unità di vita del paziente	7,5	7,5	7,5	7,5
6	8. Prelievo sangue	7,7	7,7	7,7	7,7
6	9. Rilevamento temperatura corporea	7,8	7,8	7,8	7,8
	MEDIA	7,3	7,3	7,3	7,3

Figura 10. Confronto fra i diversi aspetti da valutare relativi alle attività di infermiere UU.DD.TT. per gli anni 2019/2018/2017.



3.7 Servizio Emergenza/Urgenza 118

Il servizio in questione è l'unico, fra quelli analizzati, a registrare un lieve calo nella soddisfazione complessiva che è passata dal 9,1 del 2018 all'8,8 del 2019.

Tabella 8. *Punteggi in scala 1 – 10 attribuiti dai rispondenti ai diversi servizi erogati da Sanitaservice ASL FG relativi all'attività Emergenza/Urgenza 118. Rilevazione 2019/2018/2017.*

FREQ. ASSOL.	ATTIVITÀ DA VALUTARE EMERGENZA/URGENZA 118 anno 2019	ASPETTI DA VALUTARE			
		PROF.	COMP.	REL.	SOD.
6	1. Utilizzo di linee guida e protocolli elaborati dalla Centrale operativa per la gestione dell'emergenza/urgenza	8,8	8,7	8,8	9,0
6	2. Coordinamento con gli altri servizi pubblici adetti alle emergenze	8,8	8,8	8,7	8,7
6	3. Utilizzo degli specifici strumenti informatici, di fonìa, radiocomunicazione ed orografici	8,8	9,0	9,0	9,0
6	4. Modalità di rapporto con l'utenza, soprattutto, nei casi ad alta criticità	8,8	8,7	8,8	8,7
6	5. Livello di conoscenza della propria mansione sull'intervento	8,5	8,5	8,5	8,5
6	6. Livello di conoscenza della legge sulla privacy in relazione all'intervento	8,8	8,7	8,7	8,8
	MEDIA	8,8	8,7	8,8	8,8
FREQ. ASSOL.	ATTIVITÀ DA VALUTARE EMERGENZA/URGENZA 118 anno 2018	ASPETTI DA VALUTARE			
		PROF.	COMP.	REL.	SOD.
5	1. Utilizzo di linee guida e protocolli elaborati dalla Centrale operativa per la gestione dell'emergenza/urgenza	9,2	9,6	9,4	9,4
5	2. Coordinamento con gli altri servizi pubblici adetti alle emergenze	9,0	8,8	9,0	9,0
5	3. Utilizzo degli specifici strumenti informatici, di fonìa, radiocomunicazione ed orografici	9,0	9,0	9,0	9,0
5	4. Modalità di rapporto con l'utenza, soprattutto, nei casi ad alta criticità	9,0	9,2	9,0	9,2
5	5. Livello di conoscenza della propria mansione sull'intervento	9,6	9,2	9,2	9,2
5	6. Livello di conoscenza della legge sulla privacy in relazione all'intervento	8,6	8,8	9,0	8,8
	MEDIA	9,1	9,1	9,1	9,1

Segue Tabella 8. Punteggi in scala 1 – 10 attribuiti dai rispondenti ai diversi servizi erogati da Sanitaservice ASL FG relativi all'attività Emergenza/Urgenza 118. Rilevazione 2019/2018/2017.

FREQ. ASSOL.	ATTIVITÀ DA VALUTARE EMERGENZA/URGENZA 118 anno 2017	ASPETTI DA VALUTARE			
		PROF.	COMP.	REL.	SOD.
2	1. Utilizzo di linee guida e protocolli elaborati dalla Centrale operativa per la gestione dell'emergenza/urgenza	8,0	8,0	8,0	8,0
2	2. Coordinamento con gli altri servizi pubblici addetti alle emergenze	9,0	9,0	9,0	9,0
2	3. Utilizzo degli specifici strumenti informatici, di fonìa, radiocomunicazione ed orografici	8,0	8,0	8,0	8,0
2	4. Modalità di rapporto con l'utenza, soprattutto, nei casi ad alta criticità	8,0	8,0	8,0	8,0
2	5. Livello di conoscenza della propria mansione sull'intervento	8,5	7,5	8,0	8,0
2	6. Livello di conoscenza della legge sulla privacy in relazione all'intervento	7,5	8,5	8,0	8,0
	MEDIA	8,2	8,2	8,2	8,2

Figura 11. Confronto fra i diversi aspetti da valutare relativi alle attività di emergenza/urgenza 118 per gli anni 2019/2018/2017.



L'unico servizio che ha registrato una valutazione in calo rispetto ai due anni precedenti è quello relativo al coordinamento con gli altri servizi pubblici addetti alle emergenze.

La figura 11 evidenzia che nel 2019 i valori medi attribuiti alle 4 dimensioni valutate non sono cresciuti rispetto al 2018 ma sono aumentati rispetto al 2017.

Complessivamente, il risultato non è da considerarsi preoccupante in quanto le valutazioni medie delle 4 dimensioni giudicate sono molto elevate. È, tuttavia, opportuno analizzare con attenzione le motivazioni che hanno portato a queste valutazioni tenendo conto della delicatezza e dell'importanza di questo servizio.

4. Conclusioni

L'analisi sino ad ora fatta ci ha consentito di analizzare nel dettaglio le valutazioni ricevute da ogni servizio per le diverse dimensioni valutative utilizzate.

Nella tabella 9 abbiamo riportato un cruscotto di indicatori ottenuto calcolando i valori medi dei punteggi assegnati alle diverse attività di ciascun servizio analizzato. Nel 2019 i servizi che hanno ricevuto le valutazioni più elevate sono: Ausiliariato e Pulizia e sanificazione ambientale. Questi due servizi hanno un peso rilevante rispetto al totale degli addetti di Sanitaservice ASL FG.

Confrontando i punteggi con quelli degli anni precedenti si nota come il miglioramento sia stato graduale ma significativo con incrementi di medi di circa 0,5 punti all'anno rispetto a quasi tutte le dimensioni valutate.

L'attività di Ausiliariato UU.DD.TT. rimane quella con le performance meno soddisfacenti, visto che il punteggio medio conseguito è di 8,1 punti.

Tabella 9. *Punteggi medi pesati in scala 1 – 10 relativi ai diversi servizi erogati da Sanitaservice ASL FG negli anni di rilevazione 2019/2018/2017*

PUNTEGGI MEDI ATTRIBUITI ALLE DIVERSE ATTIVITÀ DA VALUTARE anno 2019	ASPETTI DA VALUTARE			
	PROF.	COMP.	REL.	SOD.
Servizio di ausiliariato	9,0	9,0	9,0	9,0
Servizio di pulizia e sanificazione ambientale	9,0	8,9	9,0	9,0
Servizio di manutenzione ordinaria immobili	8,6	8,4	8,5	8,5
Servizio di trasporto di farmaci e beni economici	8,8	8,8	8,8	8,8
Servizio di ausiliariato dell'UU.DD.TT.	8,1	8,1	8,0	8,1
Servizio infermieristico dell'UU.DD.TT.	8,5	8,5	8,5	8,5
Servizio di trasporto malati/feriti nella rete emergenza/urgenza (118)	8,8	8,7	8,8	8,8
MEDIA	8,8	8,7	8,8	8,8

Segue Tabella 9. *Punteggi medi pesati in scala 1 – 10 relativi ai diversi servizi erogati da Sanitaservice ASL FG negli anni di rilevazione 2019/2018/2017*

PUNTEGGI MEDI ATTRIBUITI ALLE DIVERSE ATTIVITÀ DA VALUTARE 118 anno 2018	ASPETTI DA VALUTARE			
	PROF.	COMP.	REL.	SOD.
Servizio di ausiliariato	8,5	8,5	8,5	8,5
Servizio di pulizia e sanificazione ambientale	8,3	8,2	8,3	8,1
Servizio di manutenzione ordinaria immobili	8,3	8,2	8,2	8,2
Servizio di trasporto di farmaci e beni economici	8,7	8,5	8,6	8,6
Servizio di ausiliariato dell'UU.DD.TT.	7,0	7,1	7,1	7,0
Servizio infermieristico dell'UU.DD.TT.	7,0	7,0	7,0	7,0
Servizio di trasporto malati/feriti nella rete emergenze/urgenza (118)	9,1	9,1	9,1	9,1
MEDIA	8,1	8,1	8,1	8,0
PUNTEGGI MEDI ATTRIBUITI ALLE DIVERSE ATTIVITÀ DA VALUTARE anno 2017	ASPETTI DA VALUTARE			
	PROF.	COMP.	REL.	SOD.
Servizio di ausiliariato	8,0	8,0	7,9	7,9
Servizio di pulizia e sanificazione ambientale	7,9	7,8	7,8	7,7
Servizio di manutenzione ordinaria immobili	7,4	7,5	7,4	7,3
Servizio di trasporto di farmaci e beni economici	8,1	8,0	8,1	8,1
Servizio di ausiliariato dell'UU.DD.TT.	7,0	7,0	6,9	7,0
Servizio infermieristico dell'UU.DD.TT.	7,4	7,4	7,4	7,4
Servizio di trasporto malati/feriti nella rete emergenze/urgenza (118)	8,2	8,2	8,2	8,2
MEDIA	7,7	7,7	7,7	7,7

Analoghe considerazioni possono essere fatte anche per il servizio infermieristico UU.DD.TT. che ha registrato un incremento di 1,4 punti rispetto al 2018 e di 1,1 punti rispetto al 2017. Ciò va a testimoniare il grande lavoro fatto da Sanitaservice ASL FG per migliorare i servizi prestati nelle UU.DD.TT.

Va, infine, evidenziato l'unico caso in cui i punteggi assegnati nel 2019 sono risultati inferiori rispetto a quelli del 2018. Si tratta del servizio di Emergenza/Urgenza 118, che nei confronti del 2018 perde 0,3 punti medi pur rimanendo elevato rispetto alla rilevazione effettuata nel 2017. Tale variazione non è sicuramente preoccupante soprattutto se si tiene conto che la valutazione del 2017 era pari a 8,2 e che il punteggio rilevato nel 2018 pari a 9,1 era di 0,4 punti maggiore rispetto a quello del servizio di trasporto farmaci e beni economici.

Figura 12. Confronto fra gli anni 2019/2018/2017 della soddisfazione media rispetto ai servizi erogati da Sanitaservice ASL FG.



Infine, la figura 12 fornisce, in modo molto intuitivo, informazioni sull'incremento medio registrato nei 3 anni rispetto alle 4 dimensioni valutate. L'incremento registrato nel 2019 rispetto al 2018 è quasi doppio rispetto a quello registrato fra il 2017 ed il 2018

Gli ottimi risultati ottenuti da Sanitaservice ASL FG nel 2019 saranno difficilmente replicabili in futuro visto che ormai i margini di miglioramento vanno sempre più riducendosi. In sintesi, possiamo affermare che l'analisi di customer satisfaction condotta in questi 3 anni è risultata via via più completa e dettagliata grazie all'ampiamiento della base dei rispondenti. Giova ricordare che l'indagine condotta è di tipo censuario e non campionario ed il numero delle mancate risposte è davvero poco significativo per cui i risultati ottenuti sono particolarmente robusti.

L'indagine ha evidenziato che la Sanitaservice ASL FG ha sensibilmente migliorato le proprie capacità di risolvere problemi e fornire servizi con elevati standard qualitativi ed ha dimostrato un lodevole livello di efficienza anche nella fase di raccolta dei questionari che si è svolta nel periodo di lockdown dovuto all'emergenza COVID 19.

Riferimenti bibliografici

- Hauck K., Rice N., Smith P., (2003). The influence of health care organisations on health system performance. *Journal of Health Services Research and Policy* 8(2):68–74.
- Loeb JM. (2004). The current state of performance measurement in health care. *International Journal for Quality in Health Care*, 16 (Suppl. 1):15–19.
- Mant J. (2001). Process versus outcome indicators in the assessment of quality of health care. *International Journal for Quality in Health Care*, 13(6):475–480.
- Smith P.C. (2005). Performance measurement in health care: history, challenges and prospects. *Public Money & Management*, 25(4): 213–220.
- Smith P.C. (1990). The use of performance indicators in the public sector. *Journal of the Royal Statistical Society*, 153(1): 53–72.
- Steinke C., Webster L. and Fontaine M. (2010). Evaluating building performance in healthcare facilities: an organizational perspective HERD: Health Environments. *Research & Design Journal*, 3(2): 63-83.
- Tucker M. and Smith A. (2011). User perceptions in workplace productivity and strategic FM delivery. *Facilities*, 26(5/6): 196-212.

Il carcinoma mammario in donne anziane e in donne giovani

Nunziata Ribecco^{1*}, Isabella Stasi¹, Clelia Punzo²,
Annalisa Rizzi², Giovanni Tomasicchio²

¹Dipartimento di Economia e Finanza, Università degli Studi di Bari Aldo Moro,

²Dipartimento dell’Emergenza e dei Trapianti di organi (DETO) – Bari

Riassunto: Il tema di ricerca presentato in questo articolo è stato sviluppato nell’ambito di una collaborazione con la sezione della Breast Unit del DETO ed ha riguardato l’analisi di dati del carcinoma mammario in donne anziane e in donne giovani. In particolare, si era interessati a individuare quali fattori prognostici fossero determinanti nella decisione terapeutica in entrambe le fasce d’età. In modo particolare, ci si è concentrati sulla variabile LUMINAL, che identifica da un punto di vista istochimico l’aggressività di sottotipi di carcinoma mammario. È stata studiata l’eventuale relazione tra questo fattore e le variabili dimensione del tumore (T), grado istologico del tumore (G), interessamento linfonodale (N), linfonodo sentinella (LS), TIPO INTERVENTO, PERMEAZIONE ENDOVASCOLARE PERINEURALE e METASTASI A DISTANZA. Pertanto, si sono analizzati due database di donne affette da carcinoma mammario, uno relativo a donne anziane (65-92 anni) e uno relativo a donne giovani (29-45 anni). Su entrambi i campioni, dopo un’iniziale analisi esplorativa per capire l’andamento di ciascun fattore, si è verificato se ci fosse una relazione tra alcune variabili ritenute più rilevanti con riferimento al fenomeno considerato. In conclusione, è stata condotta un’analisi multivariata, al fine di individuare quali variabili possono essere determinanti nelle scelte terapeutiche adottate.

1. Introduzione

Il carcinoma mammario rappresenta una crescita incontrollata di cellule atipiche all’interno della ghiandola mammaria. Il fenomeno del carcinoma mammario è un tema ancor oggi molto dibattuto, le cui cause non sono del tutto ben note. In generale,

* Autore corrispondente: nunziata.ribecco@uniba.it

Il lavoro qui descritto è il risultato dell’impegno congiunto degli autori. In particolare, l’analisi statistica è stata condotta da N. Ribecco e I. Stasi mentre, C. Punzo, A. Rizzi e G. Tomasicchio oltre che mettere i dati a disposizione, hanno contribuito all’interpretazione medica dei risultati.

sono stati associati alla malattia diversi fattori di rischio come: età (la maggior parte dei casi è diagnosticata in donne di età superiore a 50 anni), prima gravidanza dopo i 30 anni, menarca prima dei 12 anni, menopausa dopo i 50 anni, nulliparità e predisposizione genetica (Fondazione Umberto Veronesi – Per il Progresso delle Scienze).

In Italia, secondo l'ultima pubblicazione AIRTUM (Associazione Italiana Registri Tumori, 2019) sui numeri del cancro, il carcinoma mammario rappresenta il tumore più frequentemente diagnosticato (considerando l'intera popolazione, donne e uomini insieme). La sopravvivenza dopo la diagnosi di tumore è uno dei principali indicatori che permette di valutare l'efficacia del sistema sanitario nei confronti della patologia tumorale. La sopravvivenza, infatti, è fortemente influenzata dalla prevenzione secondaria e dalla terapia.

Pertanto, al fine di approfondire gli aspetti del fenomeno e di determinare se la variabile LUMINAL (un indicatore dei diversi sottotipi di carcinoma derivati dall'espressione congiunta di recettori ormonali e fattori di crescita) fosse relazionata alle variabili T (dimensione del tumore), G (grado istologico del tumore), N (interessamento linfonodale), LS (linfonodo sentinella), TIPO INTERVENTO, PERMEAZIONE ENDOVASCOLARE PERINEURALE e METASTASI A DISTANZA, sono stati presi in esame i dati relativi a due campioni, uno di pazienti anziane e l'altro di pazienti giovani.

2. Materiali e metodi

Il primo campione è composto da 150 pazienti, la cui età varia tra i 65 ed i 92 anni, affette da carcinoma mammario e sottoposte a trattamento chirurgico, dal 2006 al 2020, mentre il secondo campione è composto da 38 pazienti la cui età varia tra i 29 ed i 45 anni, sempre affette da carcinoma mammario, sottoposte a trattamento chirurgico dal 2006 al 2019. I dati sono stati messi a disposizione dalla prof.ssa Clelia Punzo in quanto rientrano nella sua casistica personale e riguardano donne sottoposte a trattamento chirurgico presso il reparto U.O.C. Chirurgia Bonomo (2006-2018) – U.O.S.D. Chirurgia Videolaparoscopica (2019 per le donne giovani e 2019-2020 per le donne anziane), Università degli Studi di Bari. Al fine di analizzare i dati, entrambi i database sono stati sottoposti ad alcune modifiche che si sono rese necessarie in quanto vi erano pazienti per le quali si presentavano molti dati mancanti. La soluzione adottata è stata molto accurata affinché non fossero perse informazioni importanti e non fosse compromessa la bontà dei risultati derivati dall'analisi. Per ragioni ovvie, si sono anche omessi i nomi delle pazienti.

Inizialmente sono state individuate le variabili di interesse e, ai fini dell'analisi scelta, sono state trasformate da variabili politomiche in variabili dicotomiche (sì, no; presente, assente; ecc. ...). Pertanto, la variabile PERM. ENDOVASCOLARE PERINEURALE è stata sintetizzata in due modalità SI (presenza) e NO (assenza), anche la variabile TIPO INTERVENTO è stata dicotomizzata nelle seguenti modalità: CC (chirurgia conservativa) e CD (chirurgia demolitiva). Inoltre, le modalità delle variabili LS e METASTASI A DISTANZA sono state indicate con delle etichette come di seguito specificato:

- LS con modalità CI (cellule isolate), MM (micro-metastasi), NEG (negativo), POS (positivo) e NO (non eseguito o perché forme in situ, < 2 cm, o perché forme di dimensioni maggiori e con interessamento clinico dei linfonodi);
- METASTASI A DISTANZA con modalità NO (assenti), P (polmonari), L (linfonodi), LO (linfonodi ossei), LOE (linfonodi ossei epatici), LSE (linfonodali e sospette all'encefalo), ECP (epatiche carcinosi peritoneale) e NDES (non descritte).

Le etichette delle modalità delle restanti variabili sono di seguito indicate:

- G con modalità G1 (grado istologico basso), G2 (grado istologico intermedio) e G3 (grado istologico elevato);
- T con modalità: T1 (fino a 2 cm), T1a (≤ 0.5 cm), T1b (tra 0.5 cm e 1cm), T1c (tra 1 cm e 2 cm), T2 (tra 2 cm e 5 cm), T3 (> 5 cm) e con dimensioni Tis (carcinoma in situ) e T4 che corrisponde a dimensioni molto elevate così definite: estensione diretta alla parete toracica e/o alla cute;
- N con modalità N0 (linfonodi liberi da metastasi), N1 (metastasi nei linfonodi ascellari omolaterali mobili), N1a (metastasi in 1-3 linfonodi ascellari con almeno un deposito > 0.2 cm), N1mic (micro metastasi), N2 (metastasi nei linfonodi ascellari omolaterali fissi), N2a (metastasi nei linfonodi ascellari omolaterali fissi tra di loro o ad altre strutture), N3 (metastasi in uno o più linfonodi omolaterali), N3a (metastasi nei linfonodi sottoclaveari omolaterali), NX (linfonodi non valutabili).

Considerato che lo scopo dello studio è stato quello di evidenziare quali variabili fossero determinanti nelle scelte terapeutiche, dopo un'analisi esplorativa finalizzata a conoscere l'andamento di ciascuna variabile, si è passati a verificare se esiste indipendenza fra le variabili più significative ai fini dell'analisi. Pertanto, per il campione di donne anziane, essendo la numerosità campionaria discretamente elevata, è stata effettuata una verifica d'ipotesi d'indipendenza tramite il test Chi-Quadro, mentre per il campione di donne giovani, poiché la numerosità campionaria non è molto elevata e le frequenze teoriche sono risultate spesso minori di 5, si è fatto ricorso, prevalentemente, al test esatto di Fisher (Wayne W. D., 2012).

Allorché dalla verifica d'ipotesi è risultata una significativa relazione fra le variabili sottoposte ad analisi (rifiuto dell'ipotesi nulla di indipendenza) attraverso il "corrplot", ovvero il grafico della matrice di relazione fra le variabili, è stato possibile evidenziare, per le variabili maggiormente legate tra loro, quali sono le modalità che contribuiscono a determinare, in maniera significativa, questa relazione. Le relazioni positive sono visualizzate in cerchi di colore blu ed indicano attrazione (relazione positiva), mentre le relazioni negative sono visualizzate in cerchi di colore rosso ed indicano repulsione (relazione negativa). Inoltre, l'intensità del colore e la dimensione del cerchio sono proporzionali ai contributi, ossia un colore più intenso ed un cerchio più grande evidenziano un valore elevato che è tanto più elevato quanto più aumenta l'intensità del colore e la dimensione del cerchio. Mentre, un colore meno intenso ed una circonferenza più piccola indicano un valore meno elevato che diventa sempre più piccolo al ridursi dell'intensità e della dimensione della circonferenza (Kassambara, A., 2017).

Infine, si è fatto ricorso all'Analisi delle Corrispondenze Multiple (ACM), tecnica multivariata che, vista la natura dei dati, è la più idonea a raggiungere l'obiettivo preposto. Tale tecnica, in termini generali, si propone di descrivere la struttura delle relazioni sottese alla matrice dei dati oggetto di studio attraverso la collocazione e l'analisi dei punti-modalità delle variabili in uno spazio geometrico-statistico di dimensione ridotta.

L'obiettivo è identificare sia un gruppo di individui con un profilo simile nelle loro risposte alle domande, sia le associazioni tra categorie di variabili. Nell'analisi in oggetto abbiamo scelto di focalizzare lo studio solo sulle variabili e sulle proprie modalità, poiché il nostro scopo è analizzare l'evoluzione del carcinoma mammario sulla base dei fattori che maggiormente spiegano il fenomeno, tralasciando gli individui.

Per le analisi, è stato utilizzato il software statistico R.

3. Risultati dell'analisi

3.1 Analisi su donne anziane

3.1.1 Studio delle relazioni tra le variabili

Dopo aver studiato il comportamento delle singole variabili mediante l'analisi esplorativa, siamo passati a valutare se tra le variabili ritenute più importanti ai fini dello studio del fenomeno ci fosse una qualche relazione. Pertanto, abbiamo considerato

la verifica d'ipotesi di indipendenza su tabelle di contingenza a partire dal seguente sistema d'ipotesi:

$$\begin{cases} H_0: \text{le variabili sono indipendenti} \\ H_1: \text{le variabili non sono indipendenti} \end{cases}$$

Di seguito, riportiamo i risultati delle verifiche d'ipotesi di indipendenza con riferimento alle variabili prese in esame (Tab. 1):

Tabella 1. Variabili analizzate, valore del test chi-quadro, p-value e decisione presa a seguito del risultato ottenuto.

<i>Variabili analizzate</i>	<i>Valore test χ^2</i>	<i>p-value</i>	<i>Decisione</i>
LUMINAL vs T	0.86	0.35	Accettiamo H_0 : le variabili sono indipendenti.
LUMINAL vs G	3.79	0.05	Rifiutiamo H_0: le variabili non sono indipendenti
LUMINAL vs N	0.01	0.91	Accettiamo H_0 : le variabili sono indipendenti
LUMINAL vs LS	3.79	0.05	Rifiutiamo H_0: le variabili non sono indipendenti
LUMINAL vs TIPO INTERVENTO	0.29	0.59	Accettiamo H_0 : le variabili sono indipendenti
LUMINAL vs PERM. ENDOVASCOLARE PERINEURALE	7.05	0.008	Rifiutiamo H_0: le variabili non sono indipendenti
LUMINAL vs METASTASI A DISTANZA	0.06	0.80	Accettiamo H_0 : le variabili sono indipendenti

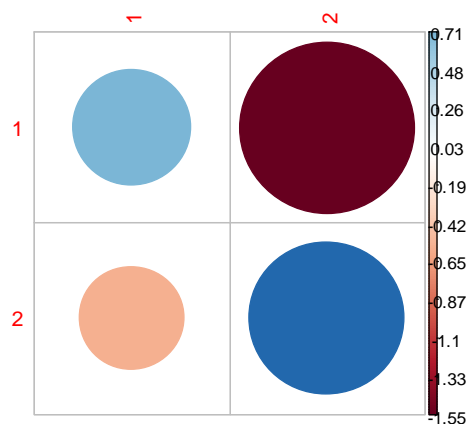
Osservando la Tab. 1, possiamo notare che le coppie di variabili risultate non significative ai fini della verifica di ipotesi, sono le seguenti: LUMINAL e T, LUMINAL ed N, LUMINAL e TIPO INTERVENTO, LUMINAL e METASTASI A DISTANZA. Al contrario, le coppie di variabili risultate significative sono: LUMINAL e G, LUMINAL e LS e LUMINAL e PERM. ENDOVASCOLARE PERINEURALE.

Pertanto, con riferimento alle coppie di variabili per cui abbiamo rifiutato l'ipotesi nulla possiamo ritenere che esiste una relazione la cui intensità potrà essere evidenziata graficamente tramite i “*corrplot*” di seguito riportati.

Come già specificato nel primo paragrafo, il “*corrplot*” è uno strumento grafico che permette di valutare in maniera immediata le relazioni fra le variabili tramite sia la colorazione delle circonferenze che la dimensione.

Analizziamo il “*corrplot*” relativo alle prime coppie di variabili risultate significativamente non indipendenti: LUMINAL e G (Fig. 1).

Figura 1. “*Corrplot*” delle variabili LUMINAL e G.

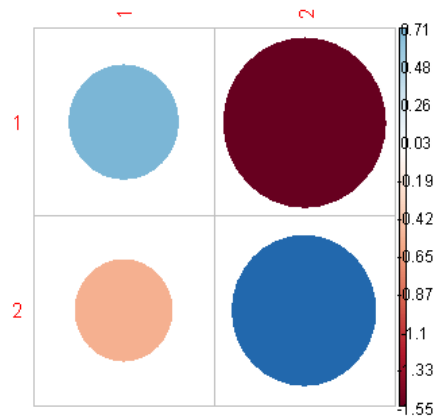


Osservando il “*corrplot*” appare immediatamente evidente che esiste un'associazione fortemente positiva tra GRADO ISTOLOGICO ELEVATO (2 in riga) e TPA (2 in colonna), seguita da una associazione sempre positiva ma meno forte tra GRADO ISTOLOGICO MEDIO-BASSO (1 in riga) e TMA (1 in colonna). Pertanto, sembrerebbe che all'aumentare dei casi in cui il grado istologico è **elevato** corrisponde un incremento di casi in cui il tumore assume forme più aggressive (TPA) mentre, all'aumentare dei casi in cui il grado istologico è **medio-basso** si osserva un aumento di casi in cui il tumore assume forme meno aggressive (TMA).

È evidente che esiste una relazione fortemente negativa tra GRADO ISTOLOGICO MEDIO-BASSO (1 in riga) e TPA (2 in colonna), seguita da una relazione sempre negativa ma meno forte tra GRADO ISTOLOGICO ELEVATO (2 in riga) e TMA (1 in colonna). Questo porta ad ipotizzare che all'aumento dei casi in cui il grado istologico è **medio-basso** corrisponde una riduzione di casi in cui il tumore assume forme più aggressive (TPA) mentre, all'aumentare dei casi in cui il grado istologico è **elevato** si presenta una riduzione dei casi in cui il tumore assume forme meno aggressive (TMA).

Con riferimento alle variabili LUMINAL e LS il “*corrplot*” seguente:

Figura 2. “*Corrplot*” delle variabili LUMINAL e LS.



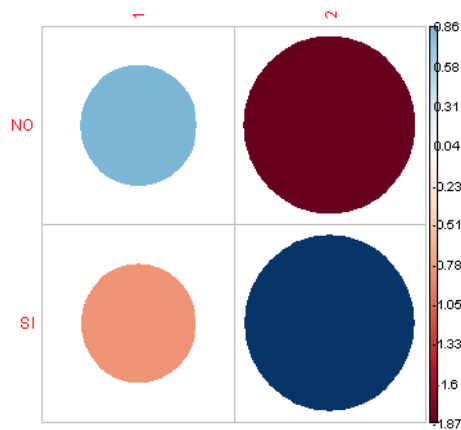
evidenzia la presenza di una relazione fortemente positiva tra LINFONODO SENTINELLA NON ESEGUITO (2 in riga) e TPA (2 in colonna), seguita da una relazione sempre positiva ma meno forte tra LINFONODO SENTINELLA ESEGUITO (1 in riga) e TMA (1 in colonna). Pertanto, sembrerebbe che all'aumentare dei casi in cui il linfonodo sentinella **non** è **eseguito**, vi sia un aumento di casi in cui il tumore assume forme più aggressive (TPA). Inoltre, all'aumento dei casi in cui il linfonodo sentinella è **eseguito**, corrisponde un aumento di casi in cui il tumore assume forme meno aggressive (TMA). Si può ritenere, quindi, che esiste una relazione fortemente negativa tra LINFONODO SENTINELLA ESEGUITO (1 in riga) e TPA (2 in colonna), seguita da una relazione sempre negativa ma meno forte tra LINFONODO SENTINELLA NON ESEGUITO (2 in riga) e TMA (1 in colonna). Possiamo concludere, quindi, che all'aumentare dei casi in cui il linfonodo sentinella è **eseguito**, corrisponde una riduzione di casi in cui il tumore assume forme più aggressive (TPA), mentre all'aumentare dei casi in cui il linfonodo sentinella **non** è **eseguito**, appare una riduzione di casi in cui il tumore assume forme meno aggressive (TMA).

Con riferimento alle variabili LUMINAL e PERM. ENDOVASCOLARE PERINEURALE appare evidente dal grafico (Fig. 3) che esiste una relazione positiva alquanto forte tra presenza di PERM. ENDOVASCOLARE PERINEURALE (SI) e TPA (2 in colonna), seguita da una relazione sempre positiva ma meno forte tra assenza DI PERM. ENDOVASCOLARE PERINEURALE (NO) e TMA (1 in colonna). Ossia, sembrerebbe che una **maggiore presenza** di PERM. ENDOVASCOLARE PERINEURALE (SI), comporti un aumento di casi in cui il tumore assume

forme più aggressive (TPA) mentre ad una **ridotta presenza** di PERM. ENDOVASCOLARE PERINEURALE (NO) corrisponde un aumento di casi in cui il tumore assume forme meno aggressive (TMA).

È evidente, inoltre, che esiste una relazione fortemente negativa tra assenza di PERM. ENDOVASCOLARE PERINEURALE (NO) e TPA (2 in colonna), seguita da una relazione sempre negativa ma meno forte tra presenza di PERM. ENDOVASCOLARE PERINEURALE (SI) e TMA (1 in colonna). Quindi si può ritenere che la presenza o meno di PERM. ENDOVASCOLARE PERINEURALE (SI/NO), comporta una riduzione di casi in cui il tumore assume forme più aggressive (TPA) ed una riduzione di casi in cui il tumore assume forme meno aggressive (TMA).

Figura 3. “Corrplot” delle variabili LUMINAL e PERM. ENDOVASCOLARE PERINEURALE



3.1.3 Analisi multivariata delle corrispondenze multiple

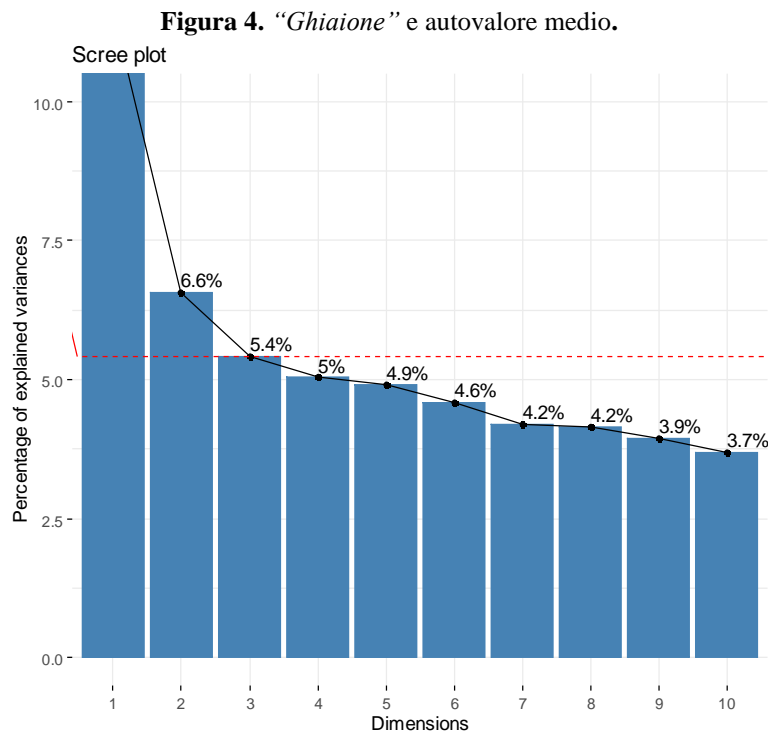
Al fine di individuare quali variabili contribuiscono maggiormente a spiegare la variabilità del fenomeno abbiamo fatto ricorso all'Analisi delle Corrispondenze Multiple. In questo paragrafo, dopo aver descritto le modalità seguite per individuare le variabili, passeremo ad analizzare i risultati ottenuti con riferimento al fenomeno oggetto di studio.

Fase 1: scelta del numero di assi da considerare

Per scegliere il numero di dimensioni da conservare per l'interpretazione dei dati non esiste una "regola empirica". Infatti, la scelta viene effettuata tramite il diagramma denominato “*ghiaione*” attraverso il quale è possibile individuare quali dimensioni contribuiscono a spiegare meglio il fenomeno. Il grafico di seguito riportato (Fig. 4)

raffigura sull'asse delle ascisse le diverse dimensioni e su quello delle ordinate la variabilità spiegata. Pertanto, le barre evidenziano la variabilità spiegata da ciascuna dimensione mentre, il punto in cui la spezzata interpolante assume la forma di una curva (il cosiddetto "*gomito*") può essere considerato come indice di una dimensionalità ottimale. Viene, inoltre, calcolato un autovalore medio che rappresenta il punto al disotto del quale i contributi per spiegare la variabilità sono considerati poco rilevanti, mentre gli assi la cui variabilità spiegata si posiziona al disopra di questo punto sono considerati importanti ai fini dell'analisi e vengono inclusi nella soluzione per l'interpretazione dei dati. (Kassambara A., 2017).

Nel grafico ("*ghiaione*") (Fig. 4) la linea rossa tratteggiata è in corrispondenza dell'autovalore medio.

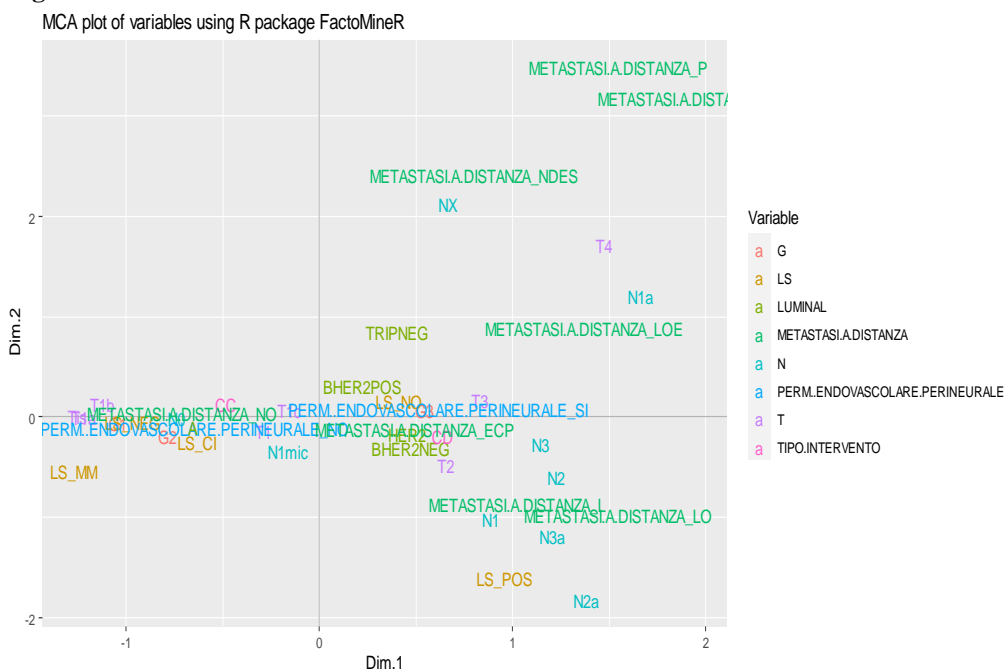


Dall'osservazione del grafico si evince che vanno prese in considerazione solo le dimensioni 1 e 2 in quanto spiegano, rispettivamente, circa l'11.6% e il 6.6% dell'inerzia totale, per un totale cumulativo del 18.2%. La dimensione 3 spiega solo il 5.4% dell'inerzia totale ed è inferiore al valore medio di equivalenza (5.41%) pertanto non viene presa in considerazione per le analisi successive.

Fase 2: coordinate delle modalità delle variabili

In questa fase si proiettano le modalità delle variabili sugli assi e si osserva la loro collocazione rispetto a questi. L'origine degli assi corrisponde alla media ponderata dei punteggi riferiti alle modalità delle variabili nella rappresentazione fattoriale che vengono definiti profili medi. I punti del grafico che si avvicinano all'origine sono quelli i cui profili sono maggiormente somiglianti con quelli medi. Analogamente i punti che presentano profili che si discostano maggiormente da quelli medi si situano in posizione periferica (Kassambara A., 2017).

Figura 5. *Proiezione sull'asse 1 e 2 delle variabili analizzate.*



La Figura 5 è il risultato della proiezione sui due assi delle variabili esaminate. In particolare possiamo evincere che, fatta eccezione per le modalità delle variabili METASTASI A DISTANZA P, METASTASI A DISTANZA LSE, METASTASI A DISTANZA NDES, NX, T4, N1a, TRIPNEG, METASTASI A DISTANZA LOE, N2a, LS POS, N3a, N1, METASTASI A DISTANZA, L e METASTASI A DISTANZA LO che sembrano isolarsi e non essere somiglianti alle altre, tutte le restanti risultano somiglianti e presentano un comportamento alquanto analogo come si può dedurre osservando il loro concentrarsi intorno all'origine degli assi.

In particolare, osserviamo che le modalità che nel grafico sono in posizioni vicine tra loro sono corrispondenti a quelle variabili che sembrano essere in relazione tra

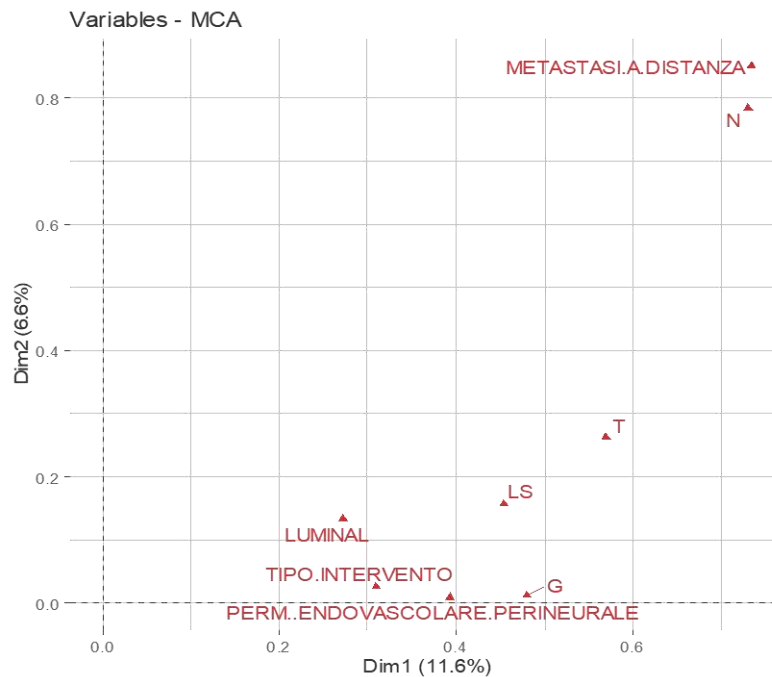
loro. Infatti, a destra dell'origine degli assi troviamo le modalità LUMINAL HER2, LUMINAL BHER2NEG, CD, METASTASI A DISTANZA ECP, T3, G3, PERM. ENDOVASCOLARE PERINEURALE SI e LS NO. Mentre, a sinistra dell'origine troviamo le modalità G2, PERM. ENDOVASCOLARE PERINEURALE NO, LS NEG, METASTASI A DISTANZA NO, N0, T1 e CC.

Le zone maggiormente concentrate inglobano le modalità posizionate nel centro del grafico. Tale rappresentazione grafica non ha l'obiettivo di determinare quali variabili contribuiscono a spiegare maggiormente la variabilità del fenomeno pertanto è opportuno proseguire con l'analisi al fine di individuare le variabili più significative ai fini dell'analisi.

Fase 3: contributo delle variabili nella definizione delle dimensioni

Una volta stabilito il numero di assi da considerare e dopo aver visualizzato le coordinate delle categorie, il riferimento prioritario va ai contributi delle variabili, cioè alla percentuale di inerzia spiegata dalla modalità della variabile, in corrispondenza di ciascun asse. L'obiettivo è determinare quali variabili contribuiscono maggiormente alla definizione delle diverse dimensioni mantenute nel modello (Kassambara A., 2017).

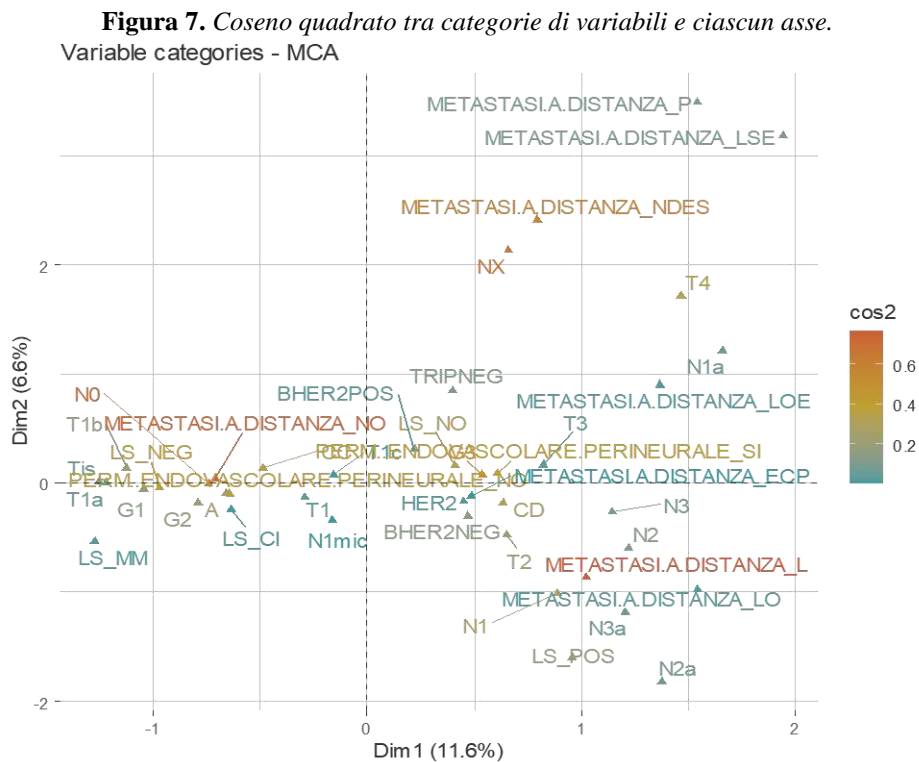
Figura 6. Variabili che contribuiscono maggiormente a definire le dimensioni del modello.



Dalla Fig. 6 si osserva che le variabili **PERM. ENDOVASCOLARE PERINEURALE**, **G** e **TIPO INTERVENTO** sono le più correlate alla dimensione 1. Al contrario, il grafico non sembra evidenziare delle variabili correlate alla dimensione 2. Inoltre, ciò che appare è che le variabili **METASTASI A DISTANZA** e **N** si discostano in maniera evidente da tutte le altre, creando una nuvola autonoma all'interno del grafico, ciò ci fa pensare che abbiano un andamento indipendente e autonomo dalle altre variabili.

Fase 4: qualità della rappresentazione di categorie variabili

La qualità della rappresentazione è chiamata “*coseno quadrato*” (**cos2**) e misura l'intensità della relazione tra categorie di variabili e un asse particolare, consentendo di individuare quale fattore descrive meglio la variabilità del fenomeno. È possibile colorare le categorie di modalità di variabili in base alla loro qualità. (Kassambara, A., 2017).



Dalla Fig. 7 osserviamo che presentano una relazione più forte con la Dimensione 1 quelle modalità che su questa direzione assumono una tonalità di colore prossimo

all'arancio ovvero N0, METASTASI A DISTANZA NO, LS NEG, PERM. ENDOVASCOLARE PERINEURALE NO, PERM. ENDOVASCOLARE PERINEURALE SI e G3. Come evidenziato in precedenza, non ci sono variabili in relazione con la Dimensione 2.

Sostanzialmente, possiamo ritenere che le modalità che meglio descrivono la variabilità del fenomeno sono:

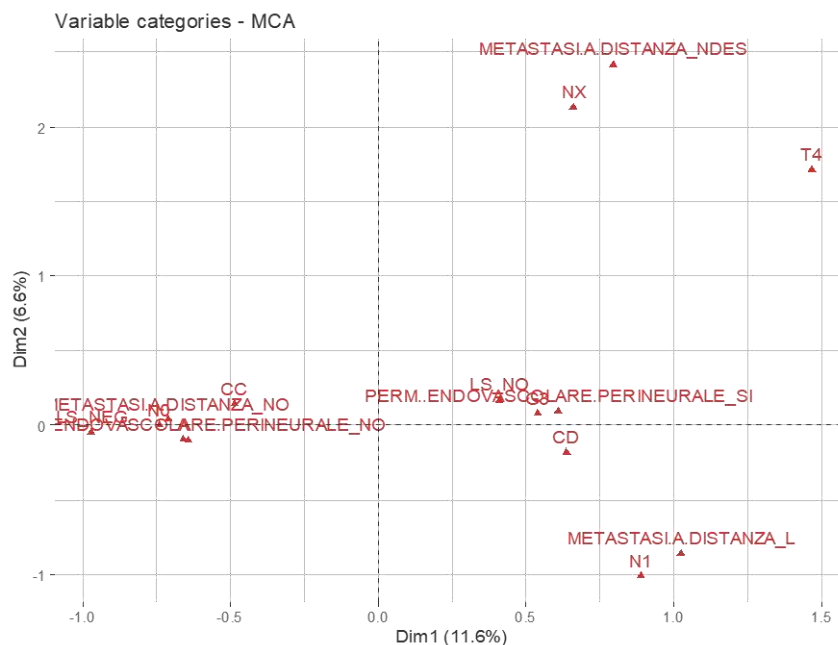
- linfonodi liberi da metastasi;
- assenza di metastasi a distanza;
- linfonodo sentinella negativo;
- assenza di permeazione endovascolare perineurale;
- presenza di permeazione endovascolare perineurale;
- grado istologico elevato.

Fase 5: scelta finale delle modalità più rappresentative

L'analisi congiunta dei contributi e dei coseni quadrati guida la scelta delle modalità da ritenere più significative ai fini della sintesi. Infatti, avendo molte categorie di variabili, è possibile visualizzarne solo alcune, le più importanti sulla base del loro indice di qualità (Kassambara, A., 2017).

Visualizziamo le prime 15 modalità di variabili con un valore più alto di cos2.

Figura 8. Modalità più rappresentative.



Nella Fig.8 ritroviamo le stesse modalità delle variabili che fino a questo momento abbiamo ritenuto fornissero un contributo importante per la spiegazione della variabilità dei dati.

In conclusione, escludendo le due modalità che non esplicitano alcun contributo (NX e METASTASI A DISTANZA NDES) in quanto modalità non valutabili, possiamo ritenere che le prime 13 modalità che contribuiscono a spiegare la variabilità del fenomeno sono:

- METASTASI A DISTANZA L (metastasi a distanza linfonodali);
- METASTASI A DISTANZA NO (assenza di metastasi a distanza);
- N1 (metastasi nei linfonodi ascellari omolaterali mobili);
- NO (linfonodi liberi da metastasi);
- LS NEG (linfonodo sentinella negativo);
- LS NO (linfonodo sentinella non eseguito);
- CC (chirurgia conservativa);
- CD (chirurgia demolitiva);
- PERM. ENDOVASCOLARE PERINEURALE SI (presenza di permeazione endovascolare perineurale);
- PERM. ENDOVASCOLARE PERINEURALE NO (assenza di permeazione endovascolare perineurale);
- G3 (grado istologico elevato);
- T4 (estensione diretta alla parete toracica e/o alla cute);
- LUMINAL A.

In modo particolare, possiamo ritenere che in presenza di permeazione endovascolare perineurale e di linfonodo sentinella non eseguito, si rileva un carcinoma mammario più aggressivo, nel quale le cellule tumorali hanno un aspetto anomalo, crescono più rapidamente e si diffondono a distanza. In questo caso, infatti, le pazienti si caratterizzano per un grado istologico molto elevato (G3) e per la chirurgia demolitiva (CD). In contrapposizione, l'assenza di metastasi a distanza, i linfonodi liberi da metastasi, il linfonodo sentinella negativo, l'assenza di permeazione endovascolare perineurale, sono indicatori di una situazione meno grave, in cui le cellule tumorali crescono lentamente e difficilmente si diffondono a distanza. In questo caso, infatti, le pazienti si caratterizzano per il luminal di tipo A e per la chirurgia conservativa (CC).

Notiamo, inoltre, che le metastasi a distanza linfonodali si posizionano vicino alla modalità N1 evidenziando una relazione.

3.2 Analisi su donne giovani

3.2.1 Studio delle relazioni tra le variabili

Anche in questo caso, dopo aver studiato il comportamento delle singole variabili mediante l'analisi esplorativa, siamo passati a valutare se tra le variabili ritenute più importanti ai fini dello studio del fenomeno ci fosse una qualche relazione. Pertanto, abbiamo considerato la verifica d'ipotesi di indipendenza su tabelle di contingenza a partire dal seguente sistema d'ipotesi:

$$\begin{cases} H_0: \text{le variabili sono indipendenti} \\ H_1: \text{le variabili non sono indipendenti} \end{cases}$$

Essendo la numerosità del campione non molto elevata ed avendo osservato in diversi casi che le frequenze teoriche erano più piccole di 5 si è fatto ricorso, per la verifica d'ipotesi di indipendenza, al test esatto di Fisher. Questa soluzione ha richiesto, inoltre, la dicotomizzazione di alcune variabili che originariamente erano poliotomiche. Tale scelta, comunque, è stata effettuata valutando attentamente che le esigenze di analisi fossero supportate da giustificazioni a carattere clinico.

Di seguito, riportiamo i risultati delle verifiche d'ipotesi di indipendenza con riferimento alle variabili prese in esame (Tab. 2):

Tabella 2. Variabili analizzate, valore del *p*-value riferito al test esatto di Fisher e decisione presa a seguito del risultato ottenuto.

<i>Variabili analizzate</i>	<i>p-value</i>	<i>Decisione</i>
LUMINAL vs T	0.20	Accettiamo H_0 : le variabili sono indipendenti
LUMINAL I vs G	0.02	Rifiutiamo H_0: le variabili non sono indipendenti
LUMINAL vs N	0.35	Accettiamo H_0 : le variabili sono indipendenti
LUMINAL vs LS	0.02	Rifiutiamo H_0: le variabili non sono indipendenti
LUMINAL vs TIPO INTERVENTO	0.17	Accettiamo H_0 : le variabili sono indipendenti
LUMINAL vs PERM. ENDOVASCOLARE PERINEURALE	0.66	Accettiamo H_0 : le variabili sono indipendenti

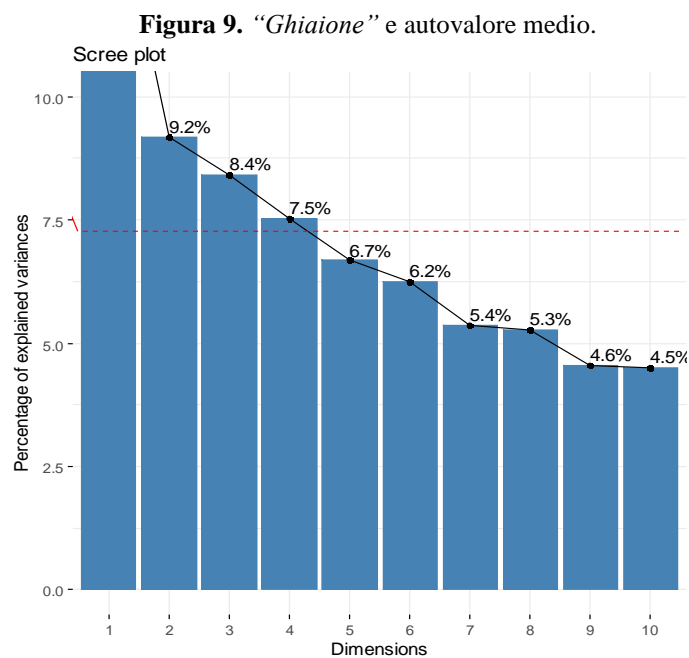
Dalla Tab. 2 si evince che le coppie di variabili risultate non significative sono le seguenti: LUMINAL e T, LUMINAL e N, LUMINAL e TIPO INTERVENTO e LUMINAL e PERM. ENDOVASCOLARE PERINEURALE. Mentre, possiamo osservare che, a differenza di quello che si è riscontrato nelle donne anziane, le uniche coppie di variabili in relazione fra loro sono le seguenti: LUMINAL e G e LUMINAL e LS.

Procedendo in maniera analoga a quanto fatto per le donne anziane, consideriamo, nel paragrafo seguente, i risultati dell'analisi multivariata.

3.2.2 Analisi multivariata delle corrispondenze multiple

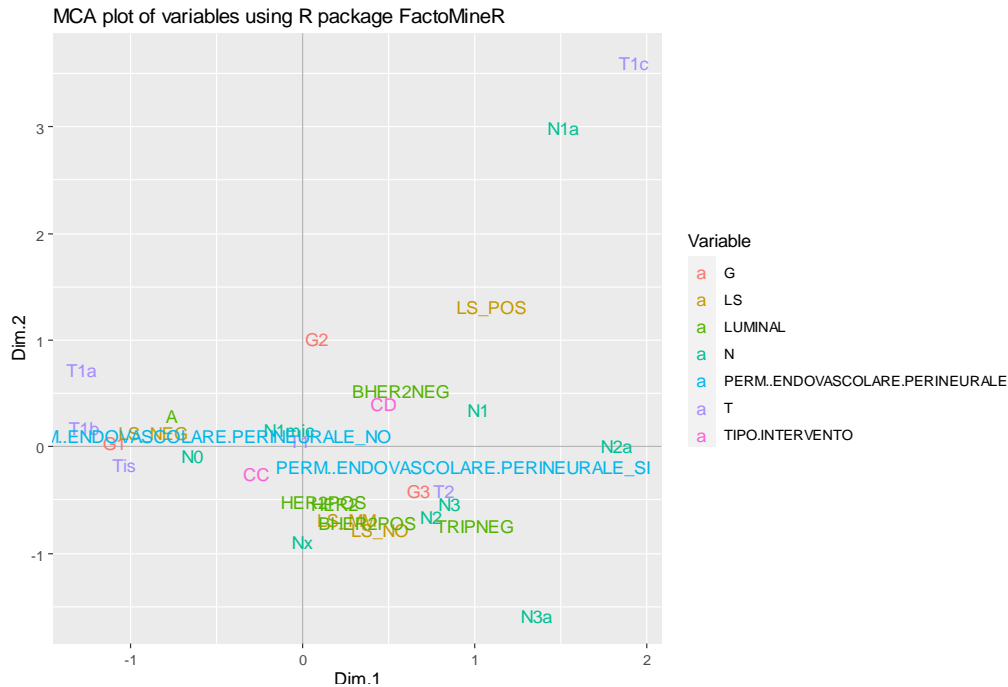
Fase 1: scelta del numero di assi da considerare

Innanzitutto, esaminiamo il diagramma “ghiaione”, in cui l'autovalore medio è in corrispondenza della linea tratteggiata rossa.



Dal grafico (Fig. 9) si evince che solo le prime quattro dimensioni spiegano, rispettivamente, circa il 14.9%, il 9.2%, l'8.4% e il 7.5% dell'inerzia totale, per un totale cumulativo del 40%. La quinta dimensione spiega il 6.7% dell'inerzia totale, inferiore al valore medio di equivalenza (7.27%), pertanto non viene considerata nelle ulteriori analisi.

Figura 10. *Proiezione sull'asse 1 e 2 delle variabili analizzate.*



Fase 2: coordinate delle modalità delle variabili

Dalla Fig. 10 possiamo evincere che, eccetto per le modalità di variabili N3a, T1c, N1a, T1a, G2 e LS POS che sembrano isolarsi e non essere somiglianti alle altre, tutte le altre tra loro hanno un comportamento molto simile, infatti si concentrano intorno all'origine degli assi.

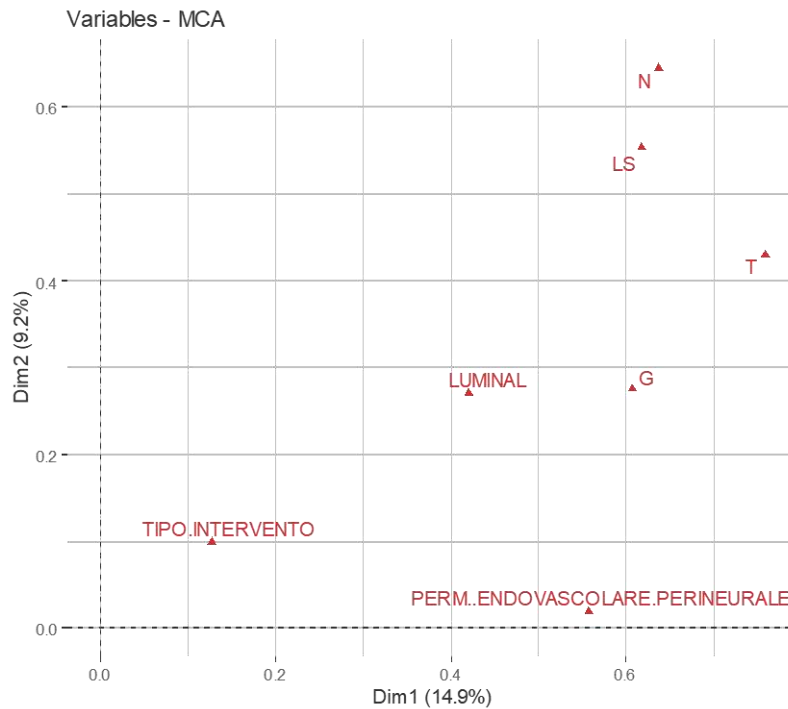
In particolare, le modalità di variabili che più sembrano avere una relazione tra loro poiché risultano vicine nel grafico sono: PERM. ENDOVASCOLARE PERINEURALE SI, G3, T2, N3, N2, LUMINAL TRIPNEG, LUMINAL HER2, LUMINAL HER2POS, LUMINAL BHER2POS e LS NO posizionate a destra rispetto all'origine. Mentre, a sinistra troviamo le modalità G1, PERM. ENDOVASCOLARE PERINEURALE NO, LS NEG, N0, N1mic, T1b, Tis, LUMINAL A e CC. Inoltre, risultano somiglianti tra loro le modalità CD e LUMINAL BHER2NEG.

Notiamo che le zone più concentrate inglobano le modalità posizionate nel centro del grafico ma non possiamo stabilire quali modalità di variabili contribuiscono maggiormente a determinare la variabilità del fenomeno in quanto questo grafico non ha questo scopo.

Fase 3: contributo delle variabili nella definizione delle dimensioni

Osserviamo i contributi delle variabili, cioè la percentuale di inerzia spiegata dalla modalità della variabile, in corrispondenza di ciascuno degli assi.

Figura 11. Variabili che definiscono maggiormente le dimensioni del modello.



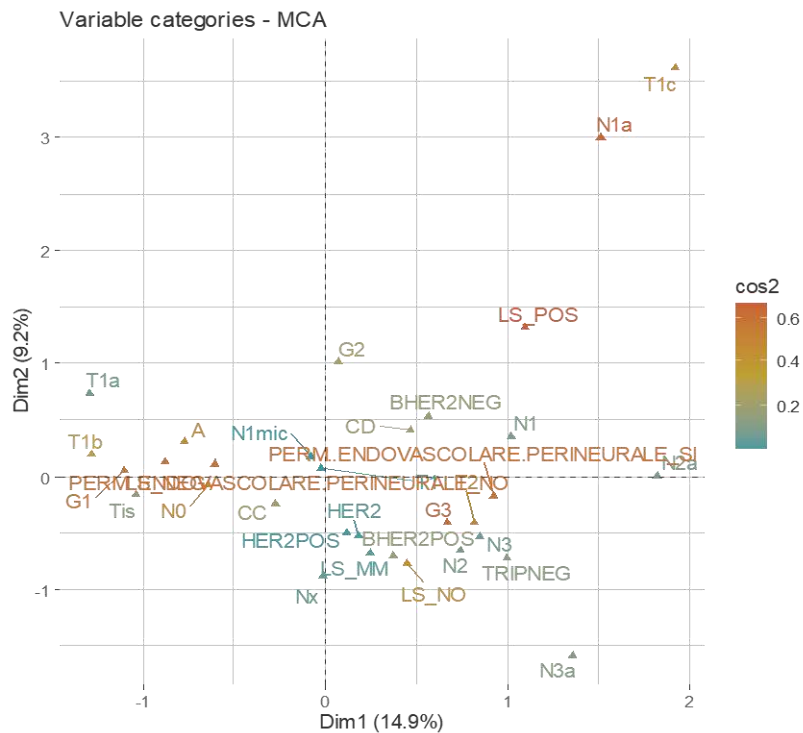
Dalla Fig. 11 si osserva che la variabile PERM. ENDOVASCOLARE PERINEURALE è la più correlata alla dimensione 1. Al contrario, il grafico non sembra evidenziare delle variabili correlate alla dimensione 2. Inoltre, ciò che appare è che le variabili N, LS e T si discostano in maniera evidente da tutte le altre, creando una nuvola autonoma all'interno del grafico e ciò ci fa pensare che abbiano un andamento indipendente e autonomo dalle altre.

Fase 4: qualità della rappresentazione di categorie variabili

La Fig. 12 evidenzia che le modalità di variabili più associate alla dimensione 1 sono quelle che su questa direzione si presentano con una tonalità di colore sull'arancio ovvero N0, LS NEG, PERM. ENDOVASCOLARE PERINEURALE NO, PERM. ENDOVASCOLARE PERINEURALE SI, G1, G3, T1b, T2 e LUMINAL A. Come

evidenziato in precedenza, non ci sono modalità di variabili associate alla dimensione 2.

Figura 12. Coseno quadrato tra categorie di variabili e ciascun asse.



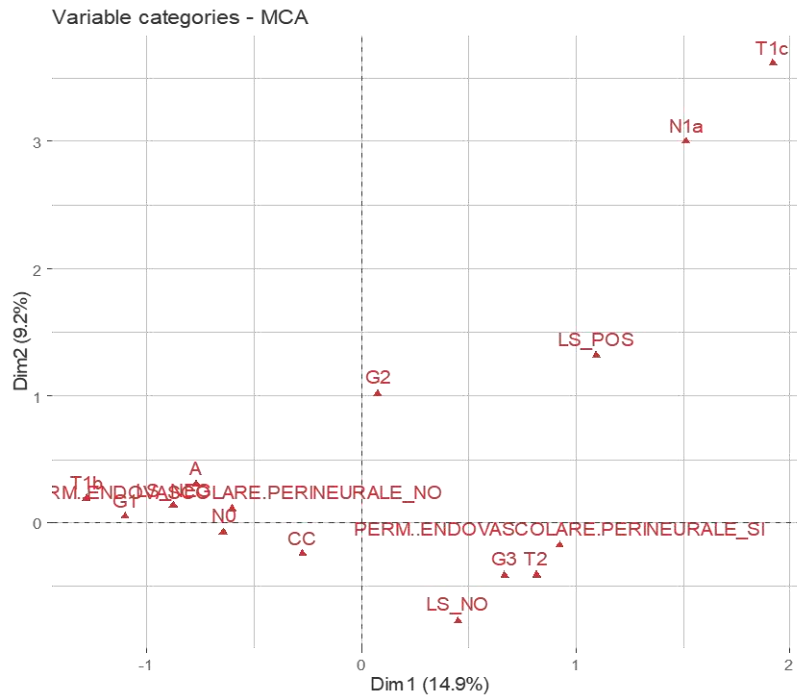
Sostanzialmente, possiamo ritenere che le modalità di variabili che meglio descrivono la variabilità del fenomeno secondo la loro qualità sono:

- linfonodi liberi da metastasi;
- linfonodo sentinella negativo;
- assenza di permeazione endovascolare perineurale;
- presenza di permeazione endovascolare perineurale;
- grado istologico basso;
- grado istologico elevato;
- dimensione del tumore compresa tra 0.5 cm e 1cm;
- dimensione del tumore compresa tra 2 cm e 5 cm nella dimensione massima;
- luminal A.

Fase 5: scelta finale delle modalità più rappresentative

Decidiamo di visualizzare le prime 15 modalità di variabili con un valore più alto di \cos^2 .

Figura 13. Modalità più rappresentative.



Nella Fig. 13 ritroviamo le stesse modalità delle variabili che fino a questo momento abbiamo ritenuto fornissero un contributo importante alla spiegazione della variabilità dei dati.

In conclusione, possiamo ritenere che le prime 15 modalità che contribuiscono a spiegare la variabilità del fenomeno sono:

- T1c (dimensione tra 1 cm e 2 cm);
- T1b (dimensione tra 0.5 cm e 1cm);
- T2 (dimensione tra 2 cm e 5 cm);
- N1a (metastasi in 1-3 linfonodi ascellari con almeno un deposito > 0.2 cm);
- N0 (linfonodi liberi da metastasi);
- G1(grado istologico basso);
- G2 (grado istologico intermedio);
- G3(grado istologico elevato);

- LS POS (linfonodo sentinella positivo);
- LS NEG (linfonodo sentinella negativo);
- LS NO (linfonodo sentinella non eseguito);
- CC (chirurgia conservativa);
- LUMINAL A;
- PERM. ENDOVASCOLARE PERINEURALE SI (presenza di permeazione endovascolare perineurale);
- PERM. ENDOVASCOLARE PERINEURALE NO (assenza di permeazione endovascolare perineurale).

In modo particolare, possiamo ritenere che in presenza di permeazione endovascolare perineurale e di linfonodo sentinella non eseguito, si rileva un carcinoma mammario più aggressivo, nel quale le cellule tumorali hanno un aspetto anomalo, crescono più rapidamente e si diffondono a distanza. In questo caso, infatti, le pazienti si caratterizzano per un grado istologico molto elevato (G3) e un tumore dalle dimensioni comprese tra 2 cm e 5 cm nella dimensione massima (T2).

In contrapposizione, l'assenza di permeazione endovascolare perineurale, i linfonodi liberi da metastasi, il linfonodo sentinella negativo e la dimensione del tumore compresa tra 0.5 cm e 1 cm sono indicatori di una situazione meno grave, in cui le cellule tumorali crescono lentamente e difficilmente si diffondono a distanza. In questo caso, infatti, le pazienti si caratterizzano per il luminal di tipo A, per la chirurgia conservativa (CC) e per un grado istologico basso (G1).

4. Conclusioni

L'analisi sin qui condotta ci porta ad alcune conclusioni riguardanti entrambi i campioni analizzati.

In particolare, con riferimento al campione di donne anziane, riteniamo che la variabile LUMINAL sia indipendente dalle variabili T, N, TIPO INTERVENTO e METASTASI A DISTANZA, mentre risulta essere in relazione con le variabili G, LS e PERM. ENDOVASCOLARE PERINEURALE. Pertanto, possiamo ritenere che in presenza di permeazione endovascolare perineurale e di linfonodo sentinella non eseguito, si rileva un carcinoma mammario più aggressivo, nel quale le cellule tumorali hanno un aspetto anomalo, crescono più rapidamente e si diffondono a distanza. In questo caso, infatti, le pazienti si caratterizzano per un grado istologico molto elevato (G3).

Nel campione da noi esaminato, una maggior percentuale era rappresentata dalle pazienti LUMINAL B. Infatti, in accordo con i dati della letteratura, nelle donne anziane, che non si sottopongono a programmi di screening, si riscontrano forme più aggressive biologicamente che vengono comunque trattate con trattamenti conservativi in rapporto alle comorbidità.

In contrapposizione, i linfonodi liberi da metastasi, il linfonodo sentinella negativo e l'assenza di permeazione endovascolare perineurale, sono indicatori di una situazione prognostica meno grave, in cui le cellule tumorali crescono lentamente e difficilmente si diffondono a distanza. In questo caso, infatti, le pazienti si caratterizzano per il LUMINAL di tipo A e rappresentano una minore percentuale.

Per quello che concerne, invece, il campione di donne giovani, riteniamo che la variabile LUMINAL sia indipendente dalle variabili T, N, TIPO INTERVENTO e PERM. ENDOVASCOLARE PERINEURALE, mentre risulta avere una relazione con le variabili G e LS. In modo particolare, possiamo ritenere che in presenza di linfonodo sentinella non eseguito, si rileva un carcinoma mammario più aggressivo, nel quale le cellule tumorali hanno un aspetto anomalo, crescono più rapidamente e si diffondono a distanza. In questo caso, infatti, le pazienti si caratterizzano per un grado istologico molto elevato (G3). In contrapposizione, i linfonodi liberi da metastasi e il linfonodo sentinella negativo sono indicatori di una situazione prognostica meno grave, in cui le cellule tumorali crescono lentamente e difficilmente si diffondono a distanza. In questo caso, infatti, le pazienti si caratterizzano per il luminal di tipo A e per un grado istologico basso (G1). È dimostrato che, a volte, pur essendo di piccole dimensioni, riscontrate tramite lo screening mammografico, sono in realtà tumori molto eterogenei, in cui prevalgono forme triple negative e ad alto grading (G3) come si evince anche dall'analisi da noi effettuata. Queste forme spesso richiedono il trattamento chemioterapico neoadiuvante prima del trattamento chirurgico.

L'analisi dei due gruppi dimostra come il carcinoma mammario non sia un'unica entità clinica e patologica e necessita di trattamenti personalizzati in base all'età. Infatti i fattori biologici ed immunoistochimici, da noi analizzati, confermano l'estrema eterogeneità del tumore mammario, che va attentamente valutata e riconosciuta.

Riferimenti bibliografici

Anastasiadi Z, Lianos GD, Ignatiadou E, Harissis H V., Mitsis M., Breast cancer in young women: an overview [Internet]. Vol. 69, *Updates in Surgery*, Springer-Verlag Italia s.r.l.; 2017 [cited 2020 Nov 1]. p. 313–7. Available from: <https://pubmed.ncbi.nlm.nih.gov/28260181/>

- Balducci L. Treatment of breast cancer in women older than 80 years is a complex task. Vol. 12, *Journal of Oncology Practice*. American Society of Clinical Oncology; 2016. p. 133–4.
- Bland M., 2009, *Statistica medica*, APOGEO, Milano.
- Ferrigni E, Bergom C, Yin Z, Szabo A, Kong AL. Breast Cancer in Women Aged 80 Years or Older: An Analysis of Treatment Patterns and Disease Outcomes [Internet]. Vol. 19, *Clinical Breast Cancer*, Elsevier Inc.; 2019 [cited 2020 Nov 1]. p. 157–64. Available from: <http://www.clinical-breast-cancer.com/article/S1526820918307055/fulltext>
- Glaser R, Marinopoulos S, Dimitrakakis C. Breast cancer treatment in women over the age of 80: A tailored approach. Vol. 110, *Maturitas*. Elsevier Ireland Ltd; 2018. p. 29–32.
- Gómez-Flores-Ramos L, Álvarez-Gómez RM, Villarreal-Garza C, Wegman-Ostrosky T, Mohar A. Breast cancer genetics in young women: What do we know? Vol. 774, *Mutation Research - Reviews in Mutation Research*. Elsevier B.V.; 2017. p. 33–45.
- Johnson HM, Irish W, Muzaffar M, Vohra NA, Wong JH. Quantifying the relationship between age at diagnosis and breast cancer-specific mortality. *Breast Cancer Res Treat*. 2019 Oct 1;177(3): 713–22.
- Kassambara A., 2017, Practical guide to principal component methods in R: PCA, M (CA), FAMD, MFA, HCPC, factoextra (Vol. 2). STHDA.
- Landis J. R., Kock G. C., 1977, The measurement of observer agreement for categorical data, *Biometrics*, 33, 159 – 174.
- Mills M, Liveringhouse C, Lee F, Nanda RH, Ahmed KA, Washington IR, et al. The prevalence of luminal B subtype is higher in older postmenopausal women with ER+/HER2- breast cancer and is associated with inferior outcomes. *J Geriatr Oncol*. 2020 Aug 26;
- Murphy BL, Day CN, Hoskin TL, Habermann EB, Boughey JC. Adolescents and Young Adults with Breast Cancer have More Aggressive Disease and Treatment Than Patients in Their Forties. *Ann Surg Oncol* [Internet]. 2019 Nov 1 [cited 2020 Nov 1];26(12):3920–30. Available from: <https://pubmed.ncbi.nlm.nih.gov/31376035/>
- Norman. G., Streiner S., 2015, *Biostatistica*, Casa Editrice Ambrosiana, Milano.
- Paluch-Shimon S, Pagani O, Partridge AH, Abulkhair O, Cardoso MJ, Dent RA, et al. ESO-ESMO 3rd international consensus guidelines for breast cancer in young women (BCY3). *Breast* [Internet]. 2017 Oct 1 [cited 2020 Nov 1];35:203–17. Available from: <https://pubmed.ncbi.nlm.nih.gov/28822332/>
- Rosner B., 2000, *Fundamentals of Biostatistics*, 5th ed. Duxbury.
- Wayne W., D., Chad L., C., 2012, *Biostatistics: A Foundation for Analysis in the Health Sciences*, 10th Edition, Casa editrice John Wiley & Sons, Inc.

Sitografia

<http://www.sthda.com/english/articles/22-principal-component-methods-videos/71-mca-in-r-using-factminer-quick-scripts-and-videos/>

<http://www.sthda.com/english/articles/31-principal-component-methods-in-r-practical-guide/114-mca-multiple-correspondence-analysis-in-r-essentials/#:~:text=MCA%20is%20generally%20used%20to,The%20associations%20between%20variable%20categories>

<http://www.sthda.com/english/articles/31-principal-component-methods-in-r-practical-guide/113-ca-correspondence-analysis-in-r-essentials/>

<https://www.fondazioneveronesi.it/magazine/tools-della-salute/glossario-delle-malattie/tumore-al-seno>

<https://www.registri-tumori.it/cms/pubblicazioni/i-numeri-del-cancro-italia-2019>

<https://www.aiom.it/linee-guida-aiom-neoplasie-della-mammella-2019/>

La povertà nei comuni del Mezzogiorno

Salvatore Cariello*, Monica Carbonara

Istat

Riassunto: Le stime sull'incidenza della povertà tra le famiglie italiane utilizzano prevalentemente i risultati di indagini campionarie, con i conseguenti limiti di significatività dei dati quando il livello di dettaglio, tematico o territoriale, si fa molto fine. Le statistiche ufficiali prodotte dall'Istat, infatti, non vanno al di là del dettaglio regionale per quanto concerne l'incidenza della povertà relativa e del rischio di povertà e si fermano al livello di ripartizione per le stime relative alla povertà assoluta. Il presente studio si propone di stimare le dimensioni del fenomeno nei comuni del Mezzogiorno utilizzando una fonte statistica, il Progetto Archivio Integrato Microdati Economici e Demografici (Arch.I.M.E.De), recentemente messa a disposizione dall'Istat e ancora poco esplorata. I data set Arch.I.M.E.De, frutto dell'integrazione di basi di dati amministrative, forniscono collezioni di microdati relativi alle caratteristiche socio-demografiche e al reddito lordo delle famiglie residenti nei comuni italiani e rappresentano, quindi, un importante strumento di conoscenza degli aspetti socio-economici dei comuni italiani.

Keywords: microdati, reddito, povertà, territorio.

1. Introduzione

“Le misure di povertà che possono essere utilizzate sono molteplici [ma] solo per alcune di esse è possibile analizzare il fenomeno a livello locale” (Regione Toscana, 2018). Il vuoto informativo delle statistiche ufficiali vincola anche la predisposizione di politiche [...] Si deduce che sia in caso di misure universalistiche sia di misure condizionate, per quantificare l'insieme dei potenziali beneficiari è necessa-

* Autore corrispondente: cariello@istat.it

rio ricorrere alle statistiche individuali di fonte amministrativa ed erariale (G. Bonanno, F. Foglia e F. Aiello, 2017).

I dati di base, infatti, sono forniti per lo più da indagini campionarie, con i conseguenti limiti di significatività quando il livello di dettaglio, tematico o territoriale, si fa molto fine. Le statistiche ufficiali prodotte dall'Istat, ad esempio, non vanno al di là del dettaglio regionale per quanto concerne l'incidenza della povertà relativa e del rischio di povertà e si fermano al livello di ripartizione per le stime relative alla povertà assoluta. A livello locale registriamo, quindi, un vuoto informativo che rende difficile, se non impossibile, la definizione di politiche territoriali e la valutazione dei loro effetti, ed è quindi giocoforza provare a utilizzare a fini statistici le fonti amministrative. Questo è il percorso che sta compiendo il team del progetto "A misura di Comune"¹ che vede coinvolti diversi attori del sistema statistico pubblico e che ha l'obiettivo di fornire un articolato set di indicatori utili per la pianificazione, programmazione e gestione degli Enti Locali. Il core delle sperimentazioni in corso è rappresentato dall'utilizzo dei microdati del Progetto Archivio Integrato Microdati Economici e Demografici (Arch.I.M.E.De) dell'Istat, insieme agli open data resi disponibili dagli Enti e dalle Amministrazioni varie. I data set Arch.I.M.E.De, frutto dell'integrazione di basi di dati amministrative, forniscono collezioni di microdati relativi alle caratteristiche socio-demografiche e al reddito lordo delle famiglie residenti nei comuni italiani.

2. Materiali e metodi

In questo lavoro, nell'ottica del calcolo di misure di contrasto della povertà, abbiamo provato ad esplorare la possibilità di produrre una misura della povertà a livello comunale utilizzando i microdati resi disponibili dal progetto Arch.I.M.E.De.

Questa nuova fonte, che utilizza appieno gli archivi amministrativi e che consente di superare i limiti di rappresentatività e dettaglio territoriale delle indagini campionarie, adoperando definizioni e classificazioni "amministrative" e non statistiche, fornisce dati non perfettamente sovrapponibili o confrontabili con quelli da indagine o non copre tutti gli ambiti delle indagini stesse. Inoltre tra le variabili re-

¹ "A misura di Comune" è sistema informativo multifonte, inserito nella sezione delle statistiche sperimentali del sito dell'Istat. Si caratterizza per essere stato progettato come un sistema di indicatori "user-oriented", selezionati e organizzati avendo come destinatari principali gli Enti locali. La sua offerta informativa è qualificata in modo fondamentale dalla componente sperimentale (Arch.I.M.E.De, Open data) e dal lavoro di ricerca, tuttora in atto, su nuove fonti e integrazione tra dati di fonti diverse.

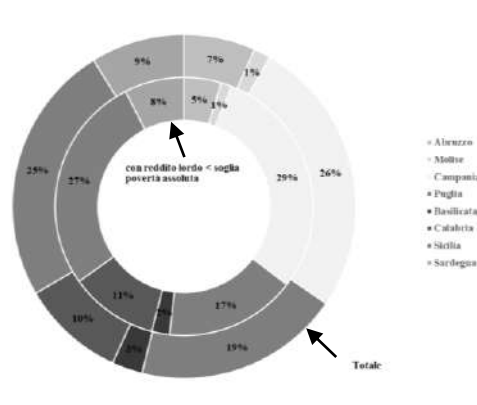
se disponibili da ArchI.M.E.De non ci sono elementi che consentano una stima dei consumi familiari. L'integrazione della Banca Dati Reddittuale del MEF, dell'Archivio Unico Persone Fisiche dell'Agenzia delle Entrate, degli archivi INPS Pensionati, Lavoratori domestici, Trattamenti monetari non pensionistici, Voucher e Uniemens, consente, invece, di ricostruire il reddito familiare, variabile alla quale ci siamo riferiti per misurare l'incidenza della povertà a livello comunale. Occorre anche sottolineare che le voci di reddito sono al lordo della tassazione e, ovviamente, non colgono il sommerso. Inoltre, il ricorso a dati amministrativi implica che le stime della povertà sono fatte nell'ambito del reddito lordo, non di quello disponibile.

Mentre nel calcolo del rischio di povertà la soglia di povertà è unica (60% del reddito mediano) e, quindi, l'attribuzione del relativo flag implica semplicemente un'operazione di confronto fra il reddito familiare e la soglia di povertà, per la povertà assoluta, invece, le soglie sono diversificate in funzione di tre fattori: la ripartizione territoriale, il tipo di comune e la composizione della famiglia per fasce d'età dei componenti. Questo rende più complessa la procedura di calcolo dell'incidenza della povertà assoluta che ha previsto:

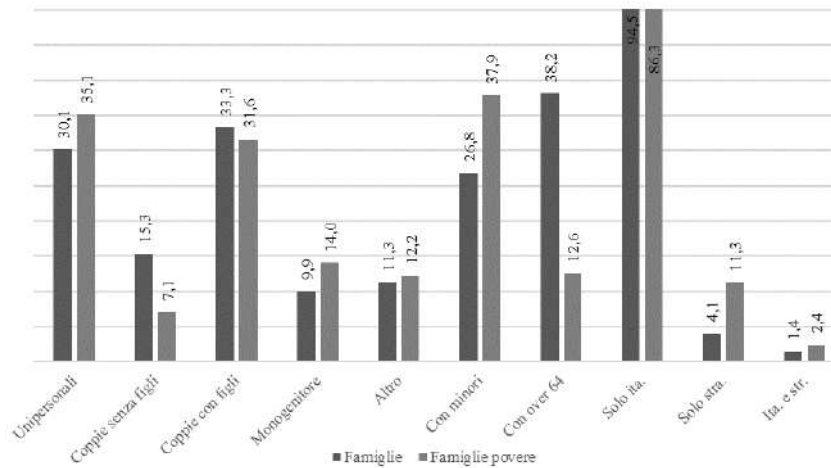
1. Classificazione dei comuni per tipologia: aree metropolitane (>250.000 abitanti), grandi comuni (>=50.000 abitanti e comuni della periferia delle aree metropolitane) e altri comuni;
2. Classificazione delle famiglie per tipologia del comune di residenza e numero di componenti per classe di età (0-3, 4-10, 11-17, 18-59, 60-74, 75 e oltre);
3. Costruzione della tabella delle soglie di povertà assoluta per i comuni del Mezzogiorno;
4. Join tra tabella delle soglie e tabella della classificazione delle famiglie e attribuzione ad ogni famiglia della corrispondente soglia di povertà;
5. Confronto fra reddito familiare lordo e soglia di povertà per l'attribuzione del relativo flag.

3. Risultati

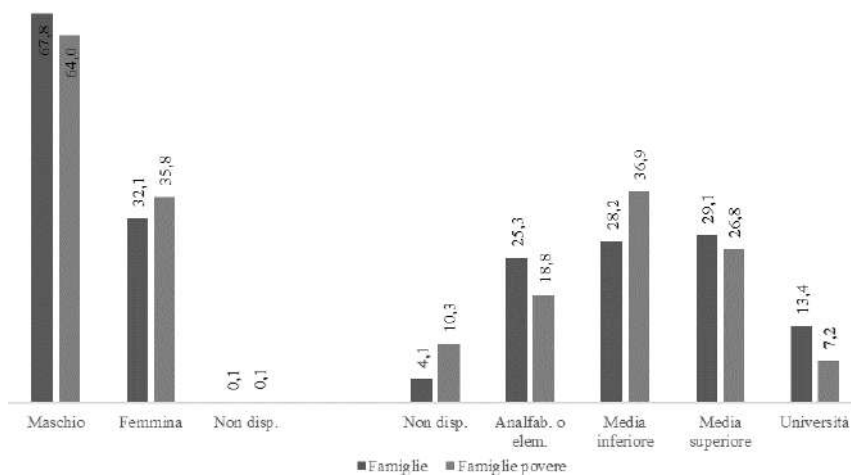
La Figura 1 evidenzia come la percentuale di famiglie povere (corona interna) sia maggiore della corrispondente quota di famiglie (corona esterna) in Campania (29% e 26%), Sicilia (27% e 25%) e Calabria (11% e 10%).

Figura 1. Famiglie e famiglie povere per regione. Anno 2016 (valori percentuali)

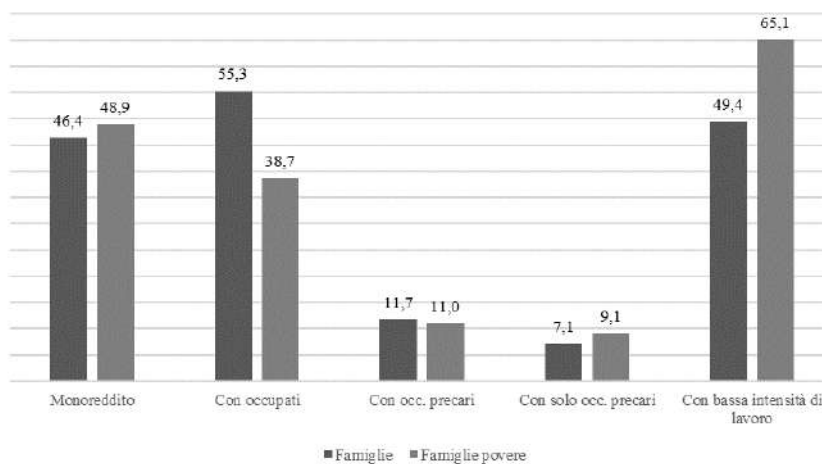
Analizzando le famiglie per le principali variabili socio-economiche (Figura 2), si evince come la povertà coinvolga principalmente le famiglie unipersonali (35,1% delle famiglie povere e il 30,1% del totale delle famiglie), quelle monogenitoriali (14% contro 9,9%), le famiglie con minori (37,9% contro 26,8) e quelle con stranieri (11,3% contro 4,1%).

Figura 2. Famiglie e famiglie povere nel Mezzogiorno per tipologia familiare e cittadinanza. Anno 2016 (valori percentuali)

La presenza di un capofamiglia femmina o con basso titolo di studio (Figura 3) concorre a determinare le condizioni di povertà familiare, come anche la presenza di una sola fonte di reddito. Le famiglie monoreddito, infatti, rappresentano il 46,4% del totale delle famiglie e il 48,9% di quelle povere (Figura 4).

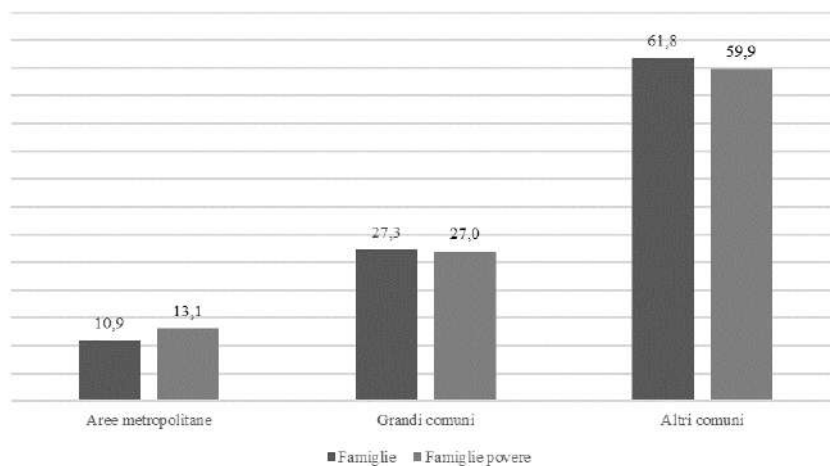
Figura 3. Famiglie e famiglie povere nel Mezzogiorno per genere e titolo di studio del capofamiglia. Anno 2016 (valori percentuali)

Le famiglie con solo occupati precari sono il 7,1% del totale delle famiglie e il 9,1% di quelle povere, quelle in cui viene utilizzato meno del 20% del potenziale di lavoro familiare sono il 49,4% del totale delle famiglie, ma ben il 65,1% di quelle povere (Figura 4).

Figura 4. Famiglie e famiglie povere nel Mezzogiorno per occupazione. Anno 2016 (valori percentuali)

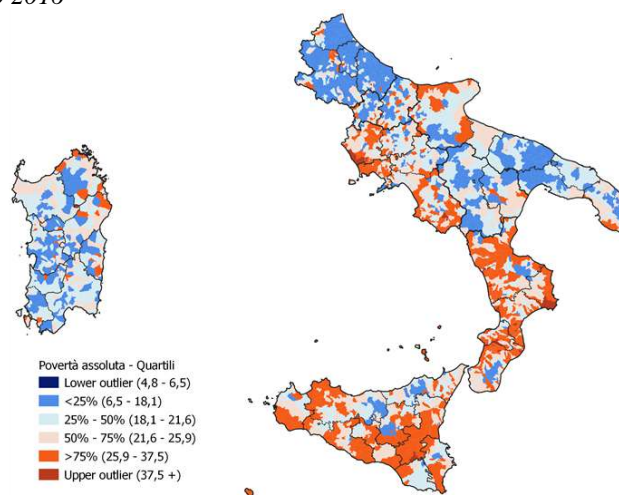
L'incidenza della povertà, infine, sembra essere più diffusa in ambito metropolitano (13,1% di famiglie povere contro 10,9% di famiglie) (Figura 5).

Figura 5. Famiglie e famiglie povere nel Mezzogiorno per tipo di comune. Anno 2016 (valori percentuali)



I comuni con maggiore incidenza della povertà assoluta (Figura 6) si concentrano soprattutto in Campania, Calabria e Sicilia. La metà dei comuni calabresi presenta una quota di famiglie povere superiore al 25,9% e in 8 comuni su 10 la percentuale di famiglie povere supera il 21,6%. In Sicilia in 4 comuni su 10 la quota di famiglie povere supera il 25,9% e in 7 su 10 supera il 21,6%. In Campania, infine, poco meno di un terzo dei comuni appartiene all'ultimo quartile (indice maggiore di 25,9) e in 6 comuni su 10 la percentuale di famiglie povere supera il 21,6%.

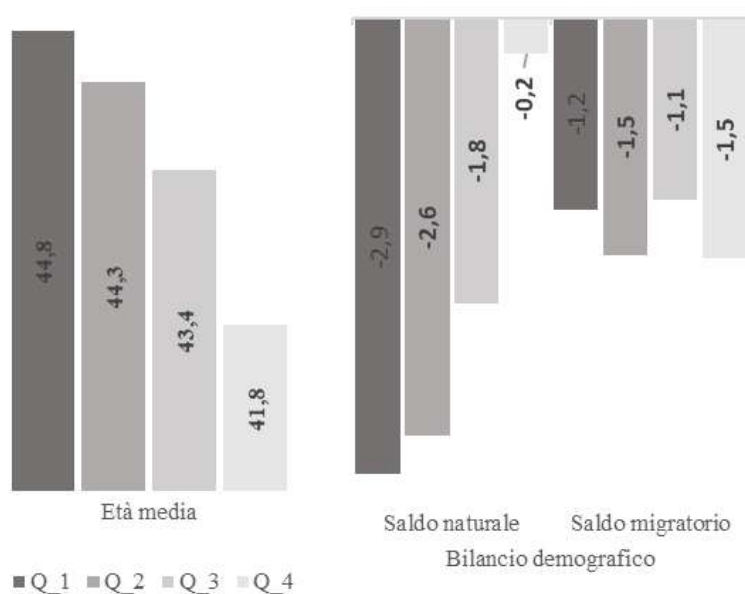
Figura 6. Percentuale di famiglie con reddito lordo inferiore alla soglia di povertà assoluta (quartili). Anno 2016



È bene ricordare che la distribuzione delle famiglie povere risulta significativamente influenzata da quella delle famiglie con reddito nullo.

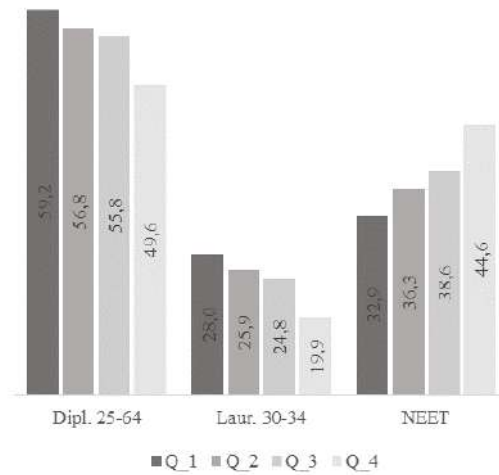
Se mettiamo in relazione l'incidenza della povertà con i principali indicatori demografici, sociali ed economici dei comuni forniti da "A misura di Comune" (Figura 7) osserviamo che i comuni con elevata incidenza di famiglie povere sono caratterizzati dalla presenza di una popolazione mediamente più giovane (l'età media dei comuni dell'ultimo quartile è pari a 41,8% contro i 44,8% dei comuni del primo quartile, le famiglie con minori sono il 26,4% contro il 24,7% del primo quartile), per una migliore stabilità della componente naturale del bilancio demografico (saldo naturale -0,2 contro il -2,9 del primo quartile) al quale si contrappone una dinamica migratoria più sostenuta (-1,5).

Figura 7. *Età media, saldo naturale e saldo migratorio per quartili dei comuni raggruppati in base all'incidenza della povertà assoluta. Anno 2016.*



I comuni con elevata incidenza di famiglie povere presentano più basse percentuali di laureati e diplomati (49,6% contro 59,2% del primo quartile) e, al contrario, elevate quote di giovani che non studiano e non lavorano (44,6% contro 32,9%) (Figura 8).

Figura 8. Percentuale di laureati, diplomati e NEET per quartili dei comuni raggruppati in base all'incidenza della povertà assoluta. Anno 2016



Per quanto concerne gli aspetti più strettamente legati al tessuto economico e al mercato del lavoro (Figura 9), i comuni in cui più di una famiglia su quattro è in condizione di povertà, si caratterizzano per un minore spirito di imprenditorialità (54,6% contro 67,1% del primo quartile) e per una struttura occupazionale meno solida (tasso di occupazione 42,8 e percentuale di famiglie a bassa intensità lavorativa al 35,6, rispetto a 55,1 e 32,8 dei comuni del primo quartile), in cui prevalgono le attività del terziario, in particolare dei servizi tradizionali (Figura 10).

Figura 9. Percentuale di occupati, di precari e di famiglie con bassa intensità di lavoro per quartili dei comuni raggruppati in base all'incidenza della povertà assoluta. Anno 2016

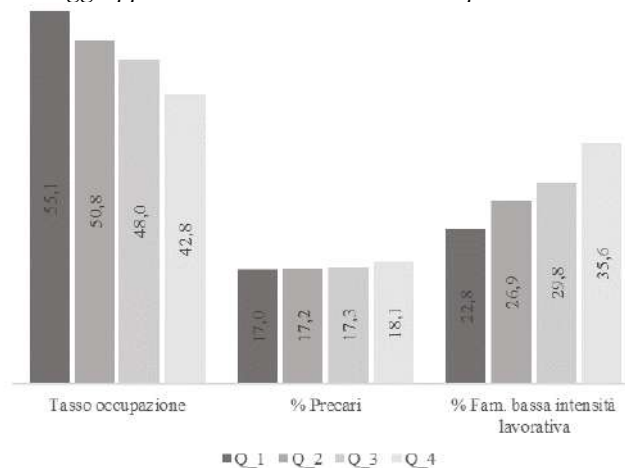
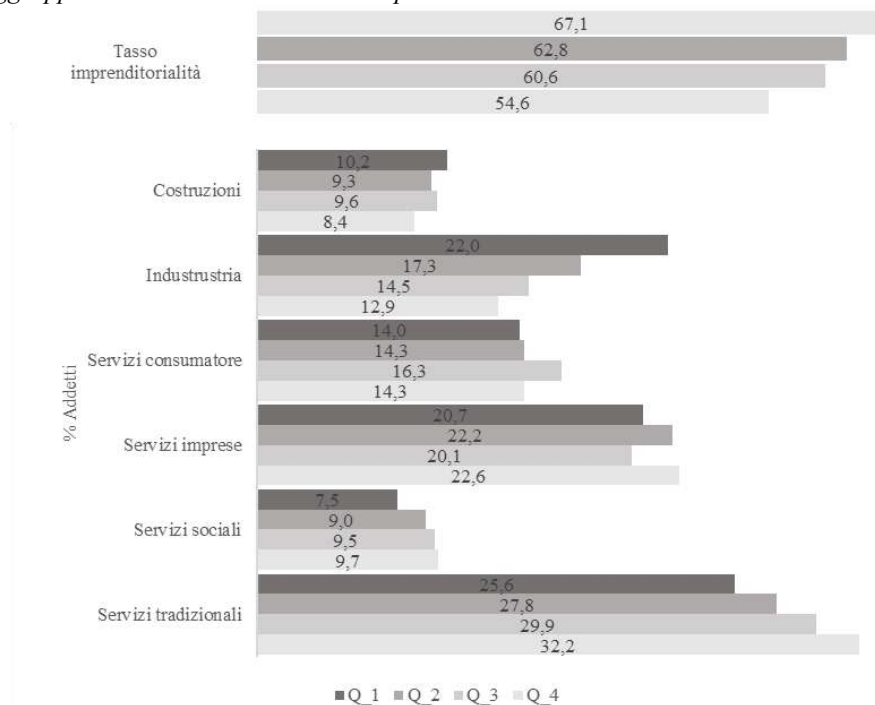
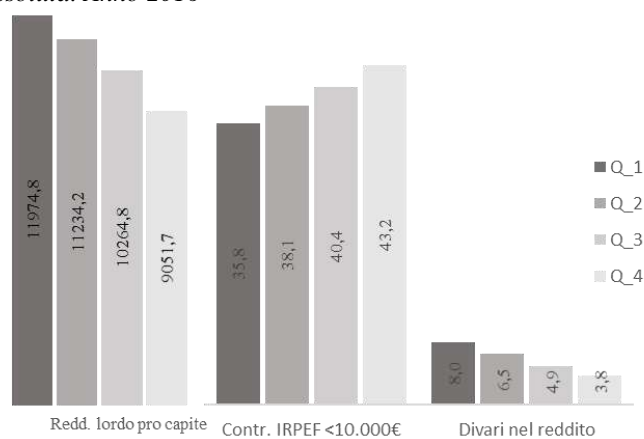


Figura 10. Tasso di imprenditorialità e struttura occupazionale per quartili dei comuni raggruppati in base all'incidenza della povertà assoluta. Anno 2016

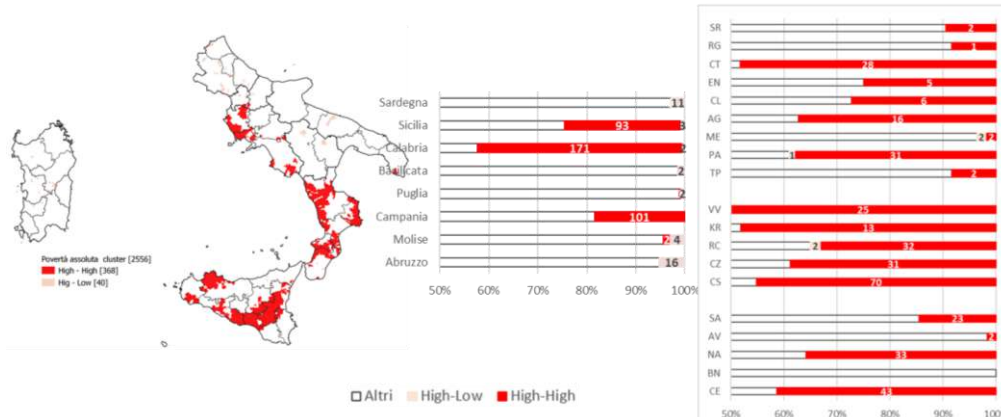
Ovviamente questi aspetti si riverberano sui livelli di reddito, infatti tra il reddito lordo pro-capite dei comuni del primo quartile e quello dei comuni del quarto c'è uno scarto di circa 3.000€ (Figura 11).

Figura 11. Reddito lordo pro-capite, percentuale di contribuenti con reddito lordo inferiore a 10.000€ e divari nel reddito per quartili dei comuni raggruppati in base all'incidenza della povertà assoluta. Anno 2016

Il cartogramma della figura 12, che riporta la distribuzione dei comuni per percentuale di famiglie povere, evidenziava la presenza di aree di particolare incidenza della povertà. Facendo ricorso alle tecniche di analisi spaziale abbiamo verificato l'esistenza di associazioni spaziali e individuato le "aree calde" della povertà nel Mezzogiorno. Il risultato dell'analisi è sintetizzato nel cartogramma della Figura 13, basato sul LISA cluster map, indicatore di associazione spaziale proposto da Anselin (1994).

Per facilitare la lettura, nel cartogramma sono riportate solo le associazioni significative dei comuni dell'ultimo quartile, colorate con una tonalità più intensa quando l'associazione è di tipo High-High, cioè comuni vicini tutti con elevato livello dell'indicatore, e con una tonalità più tenue nel caso di associazione High-Low, cioè comune con elevato livello dell'indicatore e comuni vicini a basso livello dell'indicatore.

Figura 12. Distribuzione dei comuni per cluster, regione e provincia. Anno 2016



Il cartogramma mostra chiaramente la presenza di aree di particolare incidenza del fenomeno, distribuite essenzialmente tra Campania, Calabria e Sicilia. In Campania, i 101 comuni (18,1%) che rientrano nel cluster High-High sono ubicati per la gran parte nell'area a Nord-est della regione, che dalla città Metropolitana di Napoli si spinge fino alla Pianura Campana e, più verso Ovest, verso l'area del Matese. A sud della regione troviamo 23 comuni nel Cilento-Vallo di Diano, in provincia di Salerno. In Calabria, 4 comuni su 10 rientrano nel cluster High-High. In provincia di Cosenza è presente una vasta area che dai comuni della costa tirrenica si spinge fino alla Piana di Sibari sullo Jonio, formata prevalentemente dai comuni della fascia costiera del Crotonese e, più a sud, un'ampia fascia compresa fra la

Piana di Gioia Tauro e la costa jonica. In Sicilia l'area di maggiore estensione ingloba i comuni delle province di Catania, Enna, Caltanissetta e Agrigento. Più a nord troviamo l'area della Città Metropolitana di Palermo e, in provincia di Trapani, l'area formata dai comuni di Mazara del Vallo e Castelvetro.

4. Conclusioni

I dati presentati sono le prime evidenze di un lavoro che è tuttora in corso e che si sviluppa in più direzioni. Da un lato verso la costruzione di un minimo di serie storica e verso l'estensione dello studio a tutto il territorio italiano, dall'altro verso una analisi più specifica dei casi di famiglie con reddito nullo e verso il confronto dei risultati delle elaborazioni sui microdati di ArchIMEDE con quelli di altre fonti (ad esempio la banca dati delle dichiarazioni ISEE).

Tuttavia, possiamo affermare con sicurezza che l'utilizzo a fini statistici di dati di fonte amministrativa, opportunamente integrati, offre una risposta valida ad una richiesta di informazioni statistiche sempre più dettagliate e fornisce stime di povertà a livello locale che possano supportare i "responsabili politici incaricati della pianificazione di strategie e azioni concrete nella lotta all'esclusione sociale e alla deprivazione sociale" (Giusti, Masserini, Pratesi, 2017).

Riferimenti bibliografici

- Anselin, L. (1994). Exploratory Spatial Data Analysis and Geographic Information Systems, in: M. Painho ed., *New Tools/or Spatial Analysis*, Luxembourg, Eurostat, pp. 45-54
- Bonanno, G.; Foglia, F.; Aiello, F. (2017). La povertà assoluta in Calabria. *Open Calabria*. URL <https://www.opencalabria.com/la-poverta-assoluta-calabria/>
- Cutillo, A; Raitano, M.; Siciliani, I. (2019). *Income-based and consumption-based measurement of absolute poverty: insights from Italy*. CIRET Working Papers Series - Num.2/2019
- Giusti, C.; Masserini, L.; Pratesi, M. (2017). Local Comparisons of Small Area Estimates of Poverty: An Application Within the Tuscany Region in Italy. *Social Indicators Research*, 131: 235-254
- ISTAT (2017). La povertà in Italia nel 2016. *Statistica report*

Regione Toscana (2018). *Le povertà in Toscana. Secondo Rapporto 2018*. Regione Toscana, Firenze

Rovati, G. (a cura di) (2006), *Le dimensioni della povertà*. Carocci, Roma

Sitografia

<http://amisuradicomune.istat.it>

Crisi economica, prezzi alimentari e delinquenza nella Puglia di fine Ottocento: correlazione e causalità

Ezio Ritrovato*

Università degli Studi di Bari Aldo Moro, Dipartimento di Economia e Finanza

Riassunto: Dalla crisi economica scoppiata in Puglia, dopo la riforma doganale del 1887 e la conseguente guerra commerciale con la Francia, derivarono effetti devastanti sul tenore di vita della popolazione. Ancor più in Terra di Bari, che aveva legato le sue fortune in maniera indissolubile alla produzione vinicola e alle esportazioni verso la Francia. Quando queste si bloccarono, la provincia di Bari cadde in un stato di profondo disagio economico e sociale dal quale originarono episodi di ribellione popolare e recrudescenza dei reati. Utilizzando documentazione di archivio e letteratura coeva, questo contributo intende ripercorrere gli eventi di quegli anni di fine Ottocento e spiegare eventuali nessi e causalità fra l'andamento dei prezzi agricoli e la diffusione della criminalità in Puglia e in Terra di Bari.

Keywords: Puglia; protezionismo; crisi economica; emigrazione; criminalità.

In una pubblicazione collettanea del 1900, promossa dalla Provincia di Bari in occasione dell'Esposizione Universale di Parigi, compariva un saggio dell'avvocato Michele Pantaleo intitolato *Note sulla delinquenza nelle Puglie con speciale riguardo alla provincia di Bari* (Pantaleo, 1900, 139-186). Prendendo in esame il decennio 1887-1896, l'Autore illustra l'andamento delle diverse tipologie di reati, talvolta raggruppate nelle due categorie *Reati contro la proprietà* e *Reati di sangue*, sulla scorta di dati statistici che, purtroppo, specie per la provincia di Bari

* Autore corrispondente: ezio.ritrovato@uniba.it

«non sono molti, ma sufficienti a mala pena per darci una fisionomia completa della delinquenza della nostra Terra»¹.

Una prima constatazione del carattere “regionalistico” del fenomeno porta a distinguere in Italia «due grandi zone, così per la delinquenza, come per la cultura, per la vita industriale e politica»:

L’una, la meridionale che va dal Tronto in giù e comprende l’antico regno delle Due Sicilie e lo Stato Pontificio; e l’altra che va dal Tronto in su, la settentrionale. Si può dire che man mano dalla Sicilia si sale sino alle Alpi, i foschi colori della delinquenza vanno sbiadendosi, passando dal nero rattristante delle provincie di Caltanissetta e Girgenti al chiaro rassicurante dell’alto Veneto (Pantaleo, 1900, 141)

La Puglia si presenta come un caso a parte, un’eccezione virtuosa rispetto alle altre regioni dell’Italia meridionale poiché presenta «un notevole distacco dalle regioni sorelle e quasi un punto di transizione tra il mezzogiorno e il settentrione della Penisola» (Pantaleo, 1900, 142). A supporto di tale affermazione, il Pantaleo fornisce una serie di dati che confermano per la Puglia, fra il 1890 e il 1893, una minore incidenza dei diversi reati denunciati, dalla truffa all’omicidio (Tab. 1).

Tabella 1. *Distribuzione nelle regioni meridionali dei reati denunciati tra il 1890 e il 1893. (Media annua - Incidenza su 100 mila abitanti)*

Reati	Campania e Molise	Basilicata	Abruzzi	Puglie	Calabrie	Sicilia	Sardegna	Regno
Omicidi	24,34	22,58	17,51	13,10	25,99	30,22	24,19	13,44
Lesioni	473,22	417,74	481,69	397,56	594,23	350,24	271,64	259,35
Furti	342,44	677,79	553,97	377,54	433,29	401,12	800,31	357,87
Truffe e frodi	76,06	41,27	38,37	51,06	65,99	72,03	139,99	54,27
Altri reati	2.083,93	1.581,08	1.802,83	1.393,16	2.208,97	1.446,87	2.837,02	1.442,99
Totale reati	2.999,99	2.740,46	2.894,37	2.232,42	3.328,47	2.300,48	4.073,15	2.127,92

Fonte: M. Pantaleo, 1900.

¹ Con riferimento all’approssimazione e all’inattendibilità delle indagini statistiche sulla delinquenza nelle tre province pugliesi (Capitanata, Terra di Bari e Terra d’Otranto), l’Autore lancia precise accuse alle Autorità Giudiziarie preposte a tale compito: «È doloroso constatare che i signori Procuratori Generali del Re di Trani nei loro discorsi inaugurali o non si sono occupati della distribuzione della delinquenza tra le tre provincie, o addirittura apertamente l’hanno combattuta. [...] Oh allora perché esiste una Direzione di Statistica? Perché perdere il tempo su questi lavori?» (Pantaleo, 1900, 146). Sul fenomeno della delinquenza nelle aree rurali dell’Italia e del Mezzogiorno alla fine dell’Ottocento, assumono rilevanza specifica gli studi di F. Coletti (1899 e 1910).

Quando l'analisi si sposta sulla provincia di Bari emergono le differenze con le altre province del Regno e con quelle pugliesi, ribadendo in negativo la distanza dei valori rilevati rispetto alle province del Centro-Nord ma confermando il primato di maggior diffusione della legalità nei confronti della Terra d'Otranto e della Capitanata (Tab. 2). Per spiegare «il triste primato della delinquenza» di cui si fregia la provincia di Foggia, l'Autore non riesce ad andare oltre la consueta retorica della «figlia abbandonata e reietta sulla via della civiltà, ancora smarrita nella selva selvaggia della barbarie tenebrosa» (Pantaleo, 1900, 149). Ben diverse sono, per ampiezza di illustrazione e di tematiche, le argomentazioni utilizzate quando si tratta di fornire plausibili motivazioni della condotta virtuosa delle genti di Terra di Bari. Dopo aver scomodato le teorie delle scienze etno-antropologiche e l'influenza del fattore climatico – aspetti sui quali non è il caso di soffermarsi – l'Autore giunge finalmente al «fattore economico-sociale», che è l'oggetto di questo breve contributo.

Tabella 2. *Distribuzione tra le province pugliesi dei reati denunciati tra il 1894 e il 1896. (Media annua - Incidenza su 100 mila abitanti)*

REATI	BARI	FOGGIA	LECCE
Violenze, resistenze, oltraggi all'autorità	52,84	60,97	57,95
Falsità in monete, in atti ecc.	33,92	53,13	30,61
Delitti contro il buon costume e l'ordine familiare	30,38	43,99	38,57
Omicidi di ogni specie	12,40	28,17	8,12
Lesioni personali	373,63	534,78	432,10
Furti di ogni specie	438,26	482,16	421,14
Rapine, estorsioni, ricatti	6,20	7,66	4,84
Truffe, frodi, appropriazioni indebite	71,37	70,28	63,83
Altri reati	1.239,33	1.948,81	1.578,21
Totale dei reati	2.258,33	3.229,95	2.635,37

Fonte: Pantaleo, 1900.

Ad una più generale osservazione sulla relazione inversa tra livelli di reddito pro-capite e tassi di criminalità regionali, per cui «la delinquenza discende da mezzogiorno al settentrione d'Italia, così inversamente la ricchezza discende dal nord al sud» (Pantaleo, 1900, 165), il Pantaleo fa seguire la considerazione delle «migliorate condizioni sociali degli abitanti e massime degli agricoltori» che rendono la vita in Terra di Bari «più agitata sia per lo sviluppo relativamente maggiore delle industrie e del commercio, sia per il maggiore accentramento di popolazione nei comuni» (Pantaleo, 1900, 174). Quindi, sviluppo economico e inurbamento sono

considerati fattori di contenimento delle attività criminali, al punto che, analizzando l'andamento della «dinamica della delinquenza col movimento economico della provincia di Bari nel decennio 1887-1896, non si può non concludere per l'esistenza di un rapporto di causa ed effetto».

Difatti il decennio 1887-1896 nella nostra vita agricola e commerciale presenta le stesse variabili che su per giù osservammo per la delinquenza; cioè un massimo disagio nel 1888-90, un certo rialzo di benessere nel quadriennio successivo e poi, di nuovo, disagio crescente (Pantaleo, 1900, 175).

Ma cosa aveva determinato quel «massimo disagio» nel triennio 1888-1890? Vediamo di riassumere gli eventi che, alla fine degli anni Ottanta del XIX secolo, causarono una improvvisa crisi economica in tutto il Mezzogiorno d'Italia e, in forma ancora più acuta, in Terra di Bari, per poi considerare le forme di criminalità che si manifestarono come conseguenza del profondo disagio in cui sprofondò gran parte della popolazione.

Negli anni immediatamente successivi all'introduzione della tariffa doganale del 1887² la popolazione della provincia di Bari, allora Terra di Bari, dovette affrontare privazioni e disagi dolorosissimi, in conseguenza della crisi economica che, del tutto inaspettata, si abbatté su quelle aree, devastandone l'intero tessuto economico. La narrazione che ne fa Sabino Fiorese (1900) costituisce ancora un riferimento obbligato per la conoscenza di quegli avvenimenti, delle cause che li determinarono e delle gravissime ripercussioni sui diversi settori dell'economia locale, da quello commerciale a quello creditizio. Volendo ripercorrere rapidamente i momenti più significativi del periodo, va ricordato che la chiusura del mercato francese, frutto dell'innalzamento dei dazi doganali italiani e delle conseguenti ritorsioni tariffarie da parte della Francia, aveva precluso alle esportazioni di vino barese il principale sbocco commerciale³. Nulla di irreparabile, se a quelle esportazioni non fossero state legate, con una serie di vincoli produttivi e finanziari, le sorti di gran parte dell'economia di Terra di Bari. Una dipendenza economica che si manifestava sia in agricoltura, per la poderosa opera di trasformazione culturale realizzata esclusiva-

² Ricordiamo che la nuova Tariffa generale venne approvata con la legge del 14 luglio 1887 per essere applicata dal 1° gennaio 1888, anche se effettivamente entrò in vigore dal 1° marzo di quell'anno (Stringher, 1911, 27-31).

³ Sulle vicende che condussero alla rottura dei rapporti commerciali tra Italia e Francia, vedi W. Sombart (1896, I, 245-283). Quegli stessi avvenimenti, visti da parte francese, sono descritti in S. Gerlat (1966, 273-276); Sul medesimo tema si veda anche A. Arnauné (1896).

mente in funzione dell'impianto di nuovi vigneti⁴, sia in campo creditizio, per il rapido fiorire di banche private e banche di credito cooperativo⁵ e sia, infine, nelle strategie commerciali, focalizzate sul mercato francese e del tutto lontane dalla ricerca di nuovi mercati di collocamento per l'ingente produzione vinicola⁶.

Il crollo delle esportazioni (Tab. 3), costituite prevalentemente dal vino da taglio, verso il principale mercato di sbocco, la Francia, significò pertanto la fine del periodo di prosperità vissuto nei "tre lustri d'oro" (1872-1887) e l'avvio di una crisi che avrebbe segnato profondamente e per molti anni la storia economica e sociale di Terra di Bari. La miseria e la disperazione delle classi più umili sfociarono talvolta in agitazioni popolari, in brevi tumulti e dimostrazioni di piazza, prontamente fronteggiati dalle autorità locali e dalle forze dell'ordine, sempre attente ad evitare che il malcontento e la sofferenza della gente potessero provocare sollevazioni di massa e spargimento di sangue.

Ci riferiamo agli episodi di delinquenza e di protesta popolare verificatisi in molti comuni della provincia nei primi mesi del 1889, quando gli effetti del crollo del prezzo e delle esportazioni di vino si dispiegarono in tutta la loro virulenza sulle condizioni di vita di migliaia di contadini, operai e piccoli proprietari. Le enormi quantità di vino invenduto non lasciavano speranza di reintegro di capitali investiti, di rimborso di debiti contratti, di pagamento di affitti di terre. Per alcuni significò la vendita forzata dei beni, per altri la necessità di emigrare, per molti, moltissimi, la perdita del salario di bracciante o di operaio e, quindi, la miseria e la fame.

Le prime avvisaglie dell'imminente catastrofe si presentarono già nel primo semestre del 1888, quando ormai era chiaro che non sarebbe stato più possibile vendere vino in Francia e che da ciò sarebbero derivate, nei mesi a venire, conseguenze sempre più gravi. Ne fa un elenco il Sottoprefetto di Altamura, relazionando, nel maggio 1888, al Prefetto di Bari:

⁴ In Terra di Bari, le terre coltivate a vigneto passarono dai 74.021 ettari del quinquennio 1879-1883 ai 97.371 del 1890, fino ai 140.000 del 1906, raggiungendo il 27% di tutta la superficie produttiva agraria e forestale della provincia. Cfr. Ministero di Agricoltura, Industria e Commercio (da ora M.A.I.C.) (1892).

⁵ Sul rapido sviluppo e sulla successiva crisi del sistema creditizio in provincia di Bari e in Puglia, negli anni della "corsa al vigneto", si veda M. Ottolino (2000).

⁶ Nel 1886 e nel 1887 il mercato francese assorbì il 93% e il 95% di tutte le esportazioni di vino barese, pertanto era diventato di «tale importanza e capacità di assorbimento che [la provincia di Bari] era perfettamente in grado di trascurare, o almeno non curare molto, qualsiasi altro sbocco». Cfr. A. Bertolini, E. Graziadei (1900, II, 281).

La mancanza, per molti mesi, della pioggia danneggiò le campagne, inaridì le cisterne e fé soffrire perfino la sete alle popolazioni di vari comuni del Circondario. Ma l'aggravamento della crisi si deve, senza dubbio, allo incagliato commercio vinario, che è tanta parte della vita economica di queste popolazioni. Questo incaglio ha prodotto scarsezza di capitali, ribasso dei salari a meno della metà, sospensione di lavori edilizi per conto di privati, diminuzione di lavori campestri, limitazione delle spese di famiglia, disagio e malcontento generale⁷.

Tabella 3. *Esportazioni di Terra di Bari con i principali partner commerciali (1886-'88) (valori in lire)*

	1886	1887	1888
FRANCIA			
<i>Export</i>	34.456.286	35.583.358	9.953.659
<i>(Vino)</i>	(23.670.800)	(24.440.070)	(3.828.720)
AUSTRIA-UNG.			
<i>Export</i>	6.253.240	6.463.696	6.757.178
GERMANIA			
<i>Export</i>	4.946.684	4.759.756	8.667.464
GRAN BRETAGNA			
<i>Export</i>	3.529.100	3.371.290	3.346.600
SVIZZERA			
<i>Export</i>	270.132	304.370	141.769

Fonte: Camera di Commercio e Agricoltura di Bari, 1889.

La disoccupazione bracciantile e operaia era il grande male che si diffuse in Terra di Bari dal 1888 e colpì tutti coloro che, nei comuni più popolosi o nei borghi minori, traevano dal lavoro "sotto padrone" il sostentamento per sé e per la propria famiglia. Se i proprietari e gli imprenditori si mostravano riluttanti o impossibilitati a investire in migliorie agricole o nell'edilizia urbana, l'unica speranza per la massa di disoccupati era affidata all'iniziativa pubblica. Ma i limiti dei bilanci comunali, le difficoltà e i ritardi nell'erogazione dei prestiti governativi e l'approssimarsi dell'inverno acuirono le sofferenze della popolazione che, a malapena, si sostentava con il cibo fornito dalle cucine economiche allestite da molte Amministrazioni comunali e da varie Congregazioni di carità. Per questo, l'inizio del 1889 lasciava presagire, con la prospettiva di un altro anno di crisi commerciale e agricola, il peggioramento delle condizioni di miseria in cui erano piombate le classi più umili e il riaccendersi della protesta e del malcontento che, di lì a poco, avrebbero interessato i maggiori centri della provincia.

⁷ Archivio di Stato di Bari (da ora ASB), *Fondo Intendenza-Prefettura (Agricoltura, Industria e Commercio)*, Busta 6 (*Per la crisi agraria*), Fasc. 41.

Nei primi mesi del 1889 si registrarono in diversi comuni episodi di insofferenza e di agitazione popolare. Il 13 gennaio, in un centro molto popoloso come Barletta, l'avvio di opere pubbliche che davano lavoro a circa 300 muratori, provocò una dimostrazione a favore del Pro Sindaco On. Pietrantonio Cafiero che si era adoperato per quelle iniziative. Qualche settimana più tardi si rese necessaria una riunione di proprietari e notabili – «persone tutte rispettabili per censo, per competenza della materia, per ascendente sulla popolazione e tutti uomini d'ordine» – per organizzare un pubblico comizio nel quale chiedere al Governo «la sollecita costruzione della ferrovia Barletta-Spinazzola»⁸. Una situazione abbastanza tranquilla se paragonata a quella di Andria o di Corato. Nelle campagne andriesi migliaia di contadini, ormai allo stremo, cercarono di organizzare per la fine di febbraio una grande manifestazione che avrebbe dovuto attraversare le vie del paese, chiedendo lavoro al grido di «Viva il Re, viva il Sindaco, vogliamo pane!»⁹. Gli animi si acquietarono solo quando il Consiglio Comunale, nella seduta del 2 marzo, approvò la realizzazione di alcune opere pubbliche. Purtroppo, l'inizio dei lavori, condizionato dall'erogazione di un finanziamento governativo, ritardò di alcuni giorni, e le cucine economiche allestite per fornire quasi 1.600 pasti al giorno si rivelarono insufficienti a sfamare «i circa 20.000 contadini che ora versa[va]no nella più desolante miseria»¹⁰. I rapporti delle forze dell'ordine registrano, drammaticamente, i primi morti per fame.

La miseria diventa più terribile. Il 7 andante morì per fame il contadino Monterisi, il quale lascia tre figli. La sera dell'8 il contadino Ruberti Riccardo è caduto in pubblica via sfinito per l'inedia. Fu soccorso in tempo, ed ora si spera salvarlo. Anche costui ha moglie e tre figli. Entrambe le famiglie furono prese in nota per dar loro qualche soccorso sulla cucina economica, ma simili casi, che ora accennano a ripetersi con qualche frequenza e che purtroppo si ripeteranno, esigono rimedi. Essendo Andria così popolata, la carità non può arrivare dovunque vi è il bisogno, colle risorse locali¹¹.

Andria, Corato, Barletta, Canosa, le grandi *agrotowns* pugliesi, si ritrovarono a vivere il devastante trauma dell'improvviso passaggio da uno stato di euforico e

⁸ ASB, *Fondo Intendenza-Prefettura (Agricoltura, Industria e Commercio)*, Busta 7 (*Ordine pubblico*), Fasc. 41 bis.

⁹ *Ivi*, Rapporto del Col. Caracciolo, comandante della Legione Carabinieri Reali di Bari, inviato al Prefetto il 28 febbraio, per comunicare che il pericolo della dimostrazione era stato, per il momento, scongiurato, ma «in realtà sono brutte le condizioni delle classi lavoratrici e non si può descrivere la miseria di tutti».

¹⁰ *Ivi*, Col. Caracciolo al Prefetto, 4 marzo 1889.

¹¹ *Ivi*, Col. Caracciolo al Prefetto, 10 marzo 1889.

frenetico sviluppo – con il conseguente diffondersi di un certo benessere ad ogni livello sociale – ad una depressione economica senza vie d'uscita. Si comprendono, così, anche quelle forme di sofferto orgoglio che distolgono molti, pur bisognosi, dall'accettare il pasto offerto dalle cucine economiche. È questa la situazione a Corato, dove si chiede lavoro con disperata determinazione, anche a costo di entrare nei fondi e mettersi a zappare senza esserne comandati¹².

Ma Corato sembrava, e per tale era considerata dalle autorità di Polizia, il centro di un movimento di agitazione popolare destinato a diffondersi, seppure con irrilevanti riflessi sull'ordine pubblico, in altri comuni dell'entroterra barese. Infatti, nei rapporti inviati al Prefetto di Bari, Carabinieri e Sottoprefetti riferiscono di «eccitamenti a dimostrazioni che partono da Corato, vero centro di agitazione»¹³, o della dimostrazione «promossa dal Comitato di Corato»¹⁴. Forme di protesta fin troppo moderate se si considera che provengono da un comune in cui «intere famiglie stanno digiune da più giorni e girando pei vicoli del paese si vedono scene che destano raccapriccio e che descrivendole sembrerebbe esagerazione. Non pochi, non potendo sostenersi in piedi per la fame, se ne stanno coricati nel massimo avvillimento»¹⁵.

In alcuni comuni si registrò una recrudescenza di reati, strettamente connessi ai drammatici problemi di sussistenza della popolazione. A Molfetta si verificarono ripetutamente furti ai danni di garzoni di fornai, aggrediti per strada mentre trasportavano il pane¹⁶. A Canosa il fenomeno dell'accattonaggio divenne «insistente e fino ad un certo punto allarmante»¹⁷; pratica diffusa anche ad Altamura, tanto che il Sindaco si vide costretto a far affiggere un manifesto per vietare esplicitamente tali forme di elemosina imposta:

¹² La mattina del 23 gennaio 1889, «circa 200 contadini si presentarono nel fondo del Sig. Lamonica Filippo, ad un chilometro dall'abitato, dicendo al colono di voler lavorare ad ogni costo. Questi rispose che il padrone voleva ancora attendere qualche giorno indicando i motivi, ma i contadini insistettero ed anzi 57 di essi, entrati risolutamente nei campi, cominciarono a lavorare cosicché il padrone avvertito di poi decise di lasciarli lavorare per la giornata». Il giorno dopo, sulle stesse terre, si presentarono 500 braccianti a chiedere di lavorare, ma questa volta dovettero ritirarsi di fronte all'intervento di Carabinieri, Polizia e Guardie Municipali. *Ivi*, Col. Caracciolo al Prefetto, 24 gennaio 1889.

¹³ *Ivi*, Col. Caracciolo al Prefetto, 4 marzo 1889.

¹⁴ *Ivi*, Sottoprefetto di Altamura al Prefetto, 9 febbraio 1889.

¹⁵ *Ivi*, Col. Caracciolo al Prefetto, 29 gennaio 1889.

¹⁶ *Ivi*, Col. Caracciolo al Prefetto, 9 febbraio 1889.

¹⁷ *Ivi*, Col. Caracciolo al Prefetto, 31 gennaio 1889.

Cittadini!

Il Municipio, penetrato nelle penosissime condizioni delle varie classi di operai e contadini, non ha mancato adottare tutti quei mezzi che a Lui erano consentiti per lenire, almeno in parte, tanta miseria.

Il Municipio, quindi, ha fatto quanto era in suo potere di fare e perciò è fidente che la classe bisognosa conservi quel quieto e normale contegno, che è necessario per l'ordine pubblico e di dovere per ciascun concittadino. Ove si verificasse l'incivile sistema di unirsi in gruppi per quasi coattivamente estorcere l'elemosina, massime dopo le 23 ore, allora, alle non ascoltate raccomandazioni del Magistrato Cittadino, dovrà necessariamente subentrare il rigore dell'Autorità di Pubblica Sicurezza¹⁸.

I provvedimenti adottati dalle varie Amministrazioni comunali riuscirono a lenire solo parzialmente le terribili privazioni sopportate da tanta parte della popolazione, composta nella quasi totalità da persone umili, rassegnate per atavica assuefazione ad una vita di stenti e senza speranza di riscatto. Qualche scintilla di ribellione scaturiva dal disperato bisogno, da vere e proprie esigenze di sopravvivenza, ma finiva per spegnersi in estemporanee ed effimere manifestazioni di rabbiosa impotenza.

Quando a Canosa viene recapitata al presidente della locale Banca "Principe di Napoli" una lettera anonima contenente minacce di rivolta popolare, di distruzioni e stragi qualora non fosse stato fornito lavoro ai disoccupati, il Sindaco, nel trasmetterla al Prefetto, manifesta la sua preoccupazione non tanto per una concreta possibilità di disordini, quanto per l'attività di pochi «turbolenti che, profittando delle eccezionali condizioni nelle quali si versa, cercano di gettare l'allarme e sommuovere la popolazione, che in generale è pacifica, per pescare nel torbido»¹⁹. Tuttavia, fra le espressioni veementi di una prosa semianalfabeta e dialettale, l'anonimo estensore mostra di conoscere quei meccanismi e quelle distorsioni nei rapporti produttivi che determinarono l'eccezionale gravità della crisi. A differenza della mezzadria, con il fitto a godimento o "a migliororia", l'onere delle trasformazioni colturali restava totalmente a carico dei contadini e, in caso di inadempienza nel pagamento dei canoni, il fondo ritornava al proprietario senza obbligo di indennizzo per le migliorie apportate. Nel momento della crisi, alla sofferenza e alla miseria di tanti piccoli affittuari si aggiunge anche la beffa di vedere sfumare, con l'esproprio coatto della terra, il frutto di anni di fatica e l'unica fonte di sussistenza.

Contro i grandi proprietari-usurai, contro la Banca Principe che «vuole essere pagate le cambiali» mentre «tiene tutte le nostre fatiche e a momenti si piglia i no-

¹⁸ Manifesto affisso in Altamura il 1° marzo 1889. ASB, *Fondo Intendenza-Prefettura (Agricoltura, Industria e Commercio)*, Busta 7 (*Ordine pubblico*), Fasc. 41 bis.

¹⁹ *Ivi*, Il Sindaco di Canosa al Prefetto, 4 gennaio 1889.

stri fondi», si volge il risentimento popolare e la minaccia di mettere il paese «a sacco e fuoco». Per dovere di cronaca, va riferito che della missiva incriminata venne ritenuto responsabile tale Matarrese Garibaldi, considerato individuo facinoroso e «già condannato ad un anno di carcere per violenze ad agenti della forza pubblica nella dimostrazione avvenuta a Canosa l'11 novembre 1883»²⁰. L'uso di lettere anonime ai rappresentanti del potere economico e politico locale per minacciare e, in un certo senso, preannunciare una rivolta per la cui riuscita sarebbe stato fondamentale "l'effetto sorpresa", la dice lunga sulle doti strategiche e sulle effettive intenzioni rivoluzionarie dell'anonimo mittente. È più facile percepire, fra minacce terrorizzanti, bestemmie e impropri, il senso di una supplica rivolta da chi, pur conservando quel rispetto dell'umana dignità che lo spinge a «farsi la barba a credenza per farsi vedere bello», non riesce a nascondere il retaggio di una secolare sottomissione – scrive sempre la parola "signori" con l'iniziale maiuscola – e propone l'immagine della sofferenza dei figli affamati per elemosinare quanto potrebbe pretendere in nome di giustizia e pubblica solidarietà.

Resta, comunque, indiscutibile l'incremento dei reati in Terra di Bari dopo lo scoppio della crisi e l'analisi dell'avvocato Michele Pantaleo ci porta a considerarne i collegamenti con l'andamento dei prezzi del vino e del frumento, inserendo anche il numero dei fallimenti dichiarati dalla Corte d'Appello di Trani, ai quali aggiungeremo i valori delle esportazioni complessive, sia verso altre province del Regno sia verso Paesi esteri, per costruire un prospetto riepilogativo sul quale svolgere qualche riflessione (Tab. 4).

Tabella 4. *Prezzi del vino (per hl), del grano (al ql), esportazioni, fallimenti dichiarati, reati in provincia di Bari fra il 1887 e il 1891 (valori monetari in lire) (reati = incidenza su 100 mila abitanti)*

	1887	1888	1889	1890	1891
Prezzo del vino	34,60	17,67	28,24	41,63	43,33
Prezzo del grano	22,97	22,59	23,57	24,84	26,03
Esportazioni totali	98.823.705	71.756.964	70.561.700	75.483.800	99.913.636
Fallimenti	56	118	67	46	37
Reati contro la proprietà	352,60	508,21	587,64	432,63	417,88
Reati di sangue	333,70	434,07	376,02	324,04	339,40

Fonte: Camera di Commercio e Agricoltura di Bari, 1888; M. Pantaleo, 1900, II.

²⁰ *Ivi*, Col. Caracciolo al Prefetto, 13 gennaio 1889.

Una volta acclarato il legame incontrovertibile tra tasso di criminalità e tutti gli elementi che qualificano il tenore di vita (occupazione, reddito, consumi alimentari), la prima evidente correlazione proposta dalla tabella 4 è quella tra il numero dei reati e il prezzo del vino, a sua volta determinato dall'andamento delle esportazioni. Una relazione inversa che mostra quanto le condizioni di vita della popolazione e la pace sociale fossero influenzate dalle esportazioni di vino e, di conseguenza, dal suo prezzo. Il vino era diventato la voce più importante nel movimento di esportazione dalla provincia di Bari e, pertanto, la perdita di un mercato fondamentale, come quello francese, produsse effetti devastanti su tutte le componenti di una struttura economica che fondava la propria ricchezza sull'esportazione dei prodotti agricoli. Le sorti del contadino, del fittavolo, del latifondista, del commerciante e del banchiere erano tutte legate, dopo le trasformazioni colturali degli anni Settanta e l'esplosione commerciale dei primi anni Ottanta, alla conservazione di costanti canali di smercio per il *surplus* agricolo e al sostegno delle ragioni di scambio per i prodotti primari²¹.

Lo stesso tipo di relazione inversa si ritrova con il numero dei fallimenti, manifestazione quasi immediata del crollo della principale attività economica di Terra di Bari che, come abbiamo visto, coinvolgeva imprese agricole, commerciali e creditizie. Invece per il grano, che nell'economia barese aveva un rilievo meno importante, l'aumento del prezzo influiva direttamente sul numero di reati poiché causava una riduzione dei salari reali e un impoverimento delle classi più umili, nella cui dieta il pane era la componente principale.

A questo punto, pur non essendo riferita all'andamento dei reati, appare opportuno soffermarsi su un'altra possibile relazione di causa-effetto riguardante le ripercussioni delle variazioni delle esportazioni verso i soli Paesi stranieri sui flussi migratori dalla Terra di Bari tra la fine dell'Ottocento e la Grande Guerra. Da un primo raffronto fra i dati relativi al numero degli emigrati dalla provincia di Bari tra il 1886 e il 1914 è possibile rilevare tendenze inversamente collegate fra andamento dell'emigrazione e valore delle merci esportate all'estero. Infatti, nel periodo della guerra doganale con la Francia gli emigrati passano dai 599 del 1887 ai 2.975 del 1896, mentre nello stesso decennio le esportazioni si riducono da 56.045.446 lire a 31.722.467. Nel 1897, ad una ripresa delle esportazioni fino a 50.188.172 lire,

²¹ «Il danno della crisi dai produttori si estese ai commercianti e agl'istituti di credito, e ne venne una crisi di circolazione. È noto che ogni miglioramento terriero era dovuto in parte alla fatica dei coloni, che contrassero lunghi affitti a godimento, in parte al capitale che s'era prestato di seconda e terza mano, da scontisti e da cooperative di credito, debitori della Banca Nazionale e del Banco di Napoli». Così S. La Sorsa (1915, 290).

fa seguito un calo degli espatri a 1.593 unità; ancora una caduta delle esportazioni fino al 1900 con 34.567.641 lire e un aumento degli emigrati a 2.696 unità²². A partire dal 1901, il flusso di emigrazione si fa più consistente, in concomitanza con una progressiva perdita di importanza relativa delle esportazioni sul valore complessivo del movimento commerciale ed una sempre maggiore incidenza delle importazioni²³.

Sembra possibile, quindi, individuare una relazione inversa piuttosto stabile ed evidente tra il collocamento delle produzioni baresi sui mercati esteri e i flussi migratori, che non risentono peraltro della crescita costante delle importazioni durante tutto il periodo considerato. Una spiegazione viene fornita dalla natura della quasi totalità dei prodotti inviati all'estero: si tratta di vino, olio, mandorle e derivati dalla lavorazione dell'uva e delle olive, come l'acido tartarico e l'olio al solfuro. Un anno di esportazioni in crescita era preceduto naturalmente da buoni raccolti e quindi da assorbimento di manodopera per la vendemmia o per la raccolta delle olive e delle mandorle, per l'apertura ed il funzionamento di cantine e di frantoi, per la produzione di botti da vino e fusti per olio. Non meno importante era l'impiego di forza lavoro nelle operazioni di facchinaggio e trasporto in banchina per l'imbarco sulle navi da carico. Tutto questo, contribuendo a sostenere in qualche modo il tenore di vita di gran parte della popolazione di Terra di Bari, serviva al tempo stesso ad allontanare lo spettro della partenza, spesso definitiva, per Paesi lontani.

Ma, tornando al tema centrale di questo contributo, tutte le considerazioni fin qui svolte confermano l'indissolubile legame fra il livello di benessere della popolazione e la diffusione dell'illegalità, in tutte le sue manifestazioni. Quindi, non si può non concordare con l'avvocato Pantaleo quando, concludendo il suo saggio,

²² Tutti i dati sull'emigrazione in Terra di Bari sono tratti dalle statistiche del Ministero di Agricoltura Industria e Commercio. M.A.I.C. (1891, 79); M.A.I.C. (1899, XXII); M.A.I.C. (1904, 88-93); M.A.I.C. (1906, 96-101); M.A.I.C. (1908, 100-103); M.A.I.C. (1910, 96-101); M.A.I.C. (1912, 98-103); M.A.I.C. (1914, 100-105); Ministero per l'Industria, il Commercio e il Lavoro (1918, 138-143).

²³ Negli anni 1901-1913, le importazioni sono passate, con crescita costante, da 23.532.567 a 44.381.619 lire, mentre le esportazioni, attraverso un andamento altalenante, sono passate da 38.590.542 a 37.386.865 lire; analogamente, il numero degli emigrati è aumentato o diminuito, pur con una tendenza di fondo crescente, nei periodi di minore o maggiore crescita delle esportazioni. Così se queste diminuiscono da 45.523.120 lire del 1903 a 39.833.699 nel 1906, nello stesso quadriennio gli emigrati aumentano da 8.677 a 19.414; nel 1908 le esportazioni risalgono a 49.439.679 lire e gli emigrati si riducono a 11.165; l'anno successivo le esportazioni calano a 41.780.423 lire e gli emigrati arrivano a 17.097; dal 1912 al 1913 le esportazioni diminuiscono da 46.556.706 a 37.386.865, gli emigrati aumentano da 16.494 a 26.197. Per la fonte: vedi nota precedente.

afferma che «ogni reato trova la sua spiegazione nel fattore sociale e questo fattore sociale ha la sua base fondamentale nel fattore economico. Dateci un miglioramento materiale della vita di un popolo e voi lo avrete migliorato anche intellettualmente e moralmente; dateci al contrario un peggioramento, dateci in una parola la miseria, madre della degenerazione, e voi avrete quel tale *Bouillon de culture de la criminalité* in cui il delitto si feconda, nasce e si svolge con maggiore potenzialità» (Pantaleo, 1900, 178).

Riferimenti bibliografici

- Arnauné A. (1896). La politica commerciale della Francia. *Biblioteca dell'economista*, vol. I, p. I: 219-244.
- Bertolini A., Graziadei E. (1900). La rinnovazione dei trattati e gli interessi della provincia di Bari. *Giornale degli economisti*, a. 1900, Vol. II.
- Camera di Commercio e Agricoltura di Bari (1888). *Movimento commerciale e di Navigazione della Provincia di Bari pel biennio 1886-'87*. Cannone, Bari.
- Camera di Commercio e Agricoltura di Bari (1889). *Movimento commerciale e di Navigazione della Provincia di Bari pel biennio 1888-'89*. Cannone, Bari.
- Coletti F. (1899), La delinquenza nelle classi rurali italiane secondo le più recenti statistiche pubblicate. *Bollettino quindicinale della Società degli Agricoltori Italiani*, IV, 7, 15 aprile 1899: 185-190.
- Coletti F. (1910), *Classi sociali e delinquenza in Italia nel periodo 1891-1900 con particolare considerazione delle classi rurali dell'Italia meridionale, della Sicilia e della Sardegna. Spoglio ed elaborazione nuovi di dati inediti forniti dalla Direzione generale della statistica del Regno. Introduzione e tavole statistiche con 18 cartogrammi*. Bertero, Roma.
- Fiorese S. (1900). Storia della crisi economica in Puglia dal 1887 al 1897, in AA.VV., *La Terra di Bari sotto l'aspetto storico, economico e naturale*. Vecchi, Trani, vol. II: 4-122.
- Gerlat S. (1966). Les répercussions de la rupture commerciale franco-italienne de 1887-1888: la crise économique sarde. *Cahiers d'histoire*, XI, 3: 273-276.
- La Sorsa S. (1915). *La vita di Bari durante il secolo XIX*. Vecchi, Trani, vol. II.
- Ministero di Agricoltura, Industria e Commercio (M.A.I.C.) (1891), *Statistica dell'emigrazione italiana avvenuta nell'anno 189*. Stab. Tip. dell'Opinione, Roma: 79.
- Ministero di Agricoltura, Industria e Commercio (M.A.I.C.) (1892). *Notizie e studi sulla agricoltura. Produzione e commercio del vino in Italia e all'estero*. Bertero, Roma: XXXII-XXXIII.

- Ministero di Agricoltura, Industria e Commercio (M.A.I.C.) (1899), *Statistica dell'emigrazione italiana avvenuta nell'anno 1897*. G. Bertero, Roma: XXII.
- Ministero di Agricoltura, Industria e Commercio (M.A.I.C.) (1904), *Statistica dell'emigrazione italiana per l'estero negli anni 1902 e 1903*. Bertero & C., Roma: 88-93.
- Ministero di Agricoltura, Industria e Commercio (M.A.I.C.) (1906), *Statistica dell'emigrazione italiana avvenuta negli anni 1904 e 1905*. G. Bertero & C., Roma: 96-101.
- Ministero di Agricoltura, Industria e Commercio (M.A.I.C.) (1908), *Statistica dell'emigrazione italiana avvenuta negli anni 1906 e 1907*. G. Bertero & C., Roma: 100-103.
- Ministero di Agricoltura, Industria e Commercio (M.A.I.C.) (1910) *Statistica dell'emigrazione italiana avvenuta negli anni 1908 e 1909*. G. Bertero & C., Roma: 96-101.
- Ministero di Agricoltura, Industria e Commercio (M.A.I.C.) (1912), *Statistica dell'emigrazione italiana avvenuta negli anni 1910 e 1911*. G. Bertero & C., Roma: 98-103.
- Ministero di Agricoltura, Industria e Commercio (M.A.I.C.) (1914), *Statistica dell'emigrazione italiana avvenuta negli anni 1912 e 1913*. G. Bertero & C., Roma: 100-105.
- Ministero per L'industria, il Commercio e il Lavoro (1918), *Statistica della emigrazione italiana per l'estero negli anni 1914 e 1915*. Cecchini, Roma: 138-143.
- Ottolino M. (2000). Dalle origini alla crisi dei primi anni Novanta dell'Ottocento. In M. Gangemi, M. Ottolino, M. G. Rienzo, E. Ritrovato, *La cooperazione nel credito in Puglia. Dalle origini alla Seconda guerra mondiale*. Cacucci, Bari: 3-70.
- Pantaleo M. (1900). Note sulla delinquenza nelle Puglie con speciale riguardo alla provincia di Bari, in AA.VV., *La Terra di Bari sotto l'aspetto storico, economico e naturale*. Vecchi, Trani, vol. II, Appendice: 139-186.
- Sombart W. (1896). La politica commerciale dell'Italia dall'unificazione del Regno. *Biblioteca dell'economista*, vol. I, p. I: 245-283.
- Stringher B. (1911). Gli scambi con l'estero e la politica commerciale italiana dal 1860 al 1910, in AA.VV., *Cinquant'anni di storia italiana*. Hoepli, Milano, vol. III: 27-31.

È conveniente investire nella finanza sostenibile?

Angela Maria D'Uggento^{1*}, Alessandra Milillo¹, Barbara Cafarelli²

¹Dipartimento di Economia e Finanza, Università degli Studi di Bari Aldo Moro

²Dipartimento di Economia, Management e Territorio, Università degli Studi di Foggia

Riassunto: Il presente lavoro ha lo scopo di evidenziare come, nel corso dell'ultimo decennio, sia costantemente cresciuta l'attenzione degli investitori verso i temi della sostenibilità. Gli investitori finanziari moderni ritengono un aspetto fondamentale il curare costantemente la composizione del loro portafoglio, tanto che, le molteplici opportunità di investimento offerte dai mercati sono quotidianamente monitorate allo scopo di ottenere la combinazione ideale di rischio e rendimento, in base alle caratteristiche dell'attore economico. In tale contesto, poiché le possibilità di impiego dei fondi sono divenute sempre più variegate, la gestione del rischio assume una importanza cruciale e il miglior metodo per evitare perdite superflue consiste nella diversificazione. Essendo il rischio una componente imprescindibile degli investimenti, un investitore accorto impiega i propri risparmi su più classi di attività, per sfruttare l'andamento non perfettamente concorde dei prezzi: minore è la correlazione tra i titoli, minore sarà l'effetto della rischiosità dei titoli, singolarmente considerati, sul portafoglio. Sfruttando la diversificazione, infatti, l'investitore dovrà preoccuparsi del solo rischio sistematico, cioè quello non eliminabile: in base a questo potrà dunque stimare la propria remunerazione. Le diverse tipologie di investimento si sono ampliate ulteriormente negli ultimi anni, grazie al crescente interesse nei confronti della *finanza sostenibile*. Il presente lavoro intende fornire un quadro sintetico dei punti di forza di tale settore, evidenziandone le caratteristiche peculiari con riferimento a quello degli investimenti tradizionali mediante una analisi condotta su dati a livello europeo.

Keywords: strategia di investimento; finanza sostenibile; analisi ambientale, sociale e di buon governo; redditività degli investimenti.

* Autore corrispondente: angelamaria.duggento@uniba.it. Il presente lavoro è frutto di un progetto comune, tuttavia Angela Maria D'Uggento ha provveduto alla redazione del paragrafo 1, mentre Alessandra Milillo ha redatto i paragrafi 2 e 3 e Barbara Cafarelli il paragrafo 4.

1. Introduzione

La finanza sostenibile riguarda l'applicazione del concetto di sviluppo sostenibile all'attività finanziaria. La finanza sostenibile, quindi, si pone l'obiettivo di creare valore nel lungo periodo, indirizzando i capitali verso attività che non solo generino un plusvalore economico, ma siano al contempo utili alla società e non siano a carico del sistema ambientale.

Uno dei modi in cui la sostenibilità si applica razionalmente all'attività finanziaria è la pratica dell'investimento socialmente responsabile (noto anche con l'acronimo SRI). In una accezione più ampia, gli investimenti socialmente responsabili possono essere definiti come decisioni di investimento basate su standard ambientali, sociali e di governance. Alcuni Autori (Beal et al., 2005; Valor e De la Cuesta, 2007) ritengono importante porre l'accento sulla caratteristica dell'investimento socialmente responsabile di essere anche etico ritenendo così di richiamare il processo con cui le organizzazioni attive nel sociale applicano i principi etici ad una strategia di investimento. Per altri, invece, investimento etico o fondi etici è un termine generico per descrivere gli investimenti socialmente responsabili (Kreander et al., 2005).

L'investimento sostenibile e responsabile è, quindi, una attività di gestione di strumenti finanziari orientata secondo criteri che, al perseguimento del profitto economico-finanziario, affiancano l'obiettivo di ottenere un rendimento differente, che attiene al benessere del singolo e della collettività.

L'investimento sostenibile e responsabile è definito in maniera precisa da Eurosisif, come *“un approccio agli investimenti orientato al lungo termine, che integra fattori ESG (ambientali, sociali e di governance) nel processo di ricerca, analisi e selezione di strumenti finanziari all'interno di un portafoglio di investimento. Esso combina l'analisi fondamentale con l'impegno in una valutazione dei fattori ESG allo scopo di catturare rendimenti di lungo termine per gli investitori, e di dare benefici alla società influenzando il comportamento delle aziende”*.

Alcuni studi di ricercatori indipendenti evidenziano come l'investimento responsabile e sostenibile sia in grado di produrre comunque un rendimento positivo. Occorre precisare, però, che gli investimenti responsabili sono investimenti che non hanno carattere speculativo e producono i frutti attesi nel medio-lungo termine e con rendimenti apprezzabili. Negli ultimi anni si va via via incrementando il numero delle aziende con elevati standard ambientali e di gestione tali da poter essere definite imprese sostenibili le quali, in quanto società socialmente e ambientalmente responsabili, genereranno profitto nel lungo periodo investendo in tali ambiti. A tal proposito, Desmadryl (2007) sottolinea che l'“investimento sostenibile e re-

sponsabile” è quello derivante da una strategia di investimento basata sulla performance aziendale con riferimento ai tre pilastri di sviluppo sostenibile, termine che richiama la redditività economica di lungo termine che consente alle generazioni attuali di soddisfare i loro bisogni senza mettere a rischio il soddisfacimento dei bisogni di quelle future.

Con la sottoscrizione degli Obiettivi di Sviluppo Sostenibile delle Nazioni Unite e dell'Accordo di Parigi sul clima nel 2015, l'Unione Europea ha posto la sostenibilità ambientale e sociale al centro delle proprie politiche. In particolare, l'Agenda 2030 delle Nazioni Unite, adottata dai leader mondiali nel 2015, delinea il futuro quadro di sviluppo sostenibile globale e sancisce i 17 obiettivi di sviluppo sostenibile (SDGs), tra i quali risulta prioritario il conseguimento di un pieno sviluppo sostenibile entro il 2030 a livello mondiale.

Per quanto riguarda l'Accordo di Parigi (COP21), esso ha imposto un obiettivo globale, condiviso dai 192 Paesi aderenti: ridurre le emissioni di CO2 per preservare le generazioni future e Il nuovo scenario degli investimenti ecologici si è quindi prospettato per il settore finanziario. Gli investimenti ecologici nascono con lo scopo di favorire la raccolta dei capitali per affrontare le sfide ambientali ma, al tempo stesso, rappresentano per gli investitori una differente modalità di impiego dei fondi e, quindi, una valida alternativa per la diversificazione di portafoglio. I mercati finanziari hanno accolto positivamente gli investimenti sostenibili e responsabili, i quali hanno registrato uno sviluppo esponenziale tra il 2015 ed il 2018 con una crescita prevista in ulteriore aumento per gli anni a venire.

2. Materiali e metodi

Al fine di comprendere al meglio quella che è stata nell'ultimo decennio l'evoluzione della finanza sostenibile, in particolare nel comparto del sistema bancario operante nell'area euro, saranno utilizzati i dati forniti dalla Banca Centrale Europea riguardanti l'intero collettivo delle banche europee, costituito all'incirca da 4.500 aziende, e i rapporti annuali editi dalla Fondazione Finanza etica. All'interno del suddetto collettivo delle banche europee, solo 23 banche sono state classificate come etiche e sostenibili. L'individuazione di queste ultime è avvenuta tra quelle che hanno reso disponibili i bilanci di almeno sette degli ultimi dieci anni e che mostrassero un prevalente orientamento sociale e ambientale.

Inoltre, data la notevole mole di dati e le differenti modalità di rilascio, per esigenze computazionali, si è proceduto sulla base di due assunzioni preliminari: as-

sumere solo gli ultimi tre lustri come periodi di riferimento nell'analisi temporale delle performance e, per esigenze di confrontabilità tra i dati dei differenti anni, si è seguita la metodologia di Finanza etica che ha optato per cambiare la precedente impostazione di analisi, basata sul confronto tra le banche etiche e sostenibili europee e le grandi banche “sistemiche¹”, nella comparazione tra l'andamento del sistema bancario europeo nel suo complesso con quello delle sole banche etiche. La ragione di tale scelta risiede nell'aver constatato in passato delle differenze “strutturali e non soggette a cambiamenti significativi” nel medio periodo. Inoltre, come chiarito nel rapporto 2020, confrontare le piccole banche etiche con i giganti del sistema bancario europeo implicherebbe comparare soggetti profondamente diversi per dimensione, ragione sociale e struttura dell'azionariato. Gli autori (Cavallito et al, 2020) ritengono più utile un confronto “meno polarizzato” e sostituiscono il dato delle banche sistemiche con l'intero aggregato delle banche dell'area euro, di cui, tuttavia, le banche etiche fanno parte.

Lo scopo della presente analisi è stato quello di capire se le banche etiche e sostenibili, che finanziano progetti sociali, ambientali e culturali, siano anche solide sotto il profilo finanziario-patrimoniale e dotate di redditività sotto il profilo economico e, conseguentemente, stabilire se registrino prestazioni simili o diverse da quelle delle banche “tradizionali”.

3. Risultati

Il primo aspetto analizzato è stato il peso dell'attività creditizia sul totale delle attività, confrontando, come precisato, tale dato per le banche etiche e sostenibili e per l'aggregato “banche europee”, corrispondente al sistema bancario europeo nel suo complesso.

Tabella 1. *Peso percentuale-delle attività creditizie sul totale delle attività rispetto alla tipologia di banche (etiche vs complesso)*

Tipologia di banche rispetto alla caratteristica della eticità	Crediti/Totale attivo		
	2008	2013	2018
Banche etiche/sostenibili europee	72,28	73,54	76,11
Banche europee	37,12	38,70	39,80

Fonte: Fondazione Banca Etica

¹ Per banche sistemiche si intendono le grandi banche che, per dimensioni, potrebbero causare gravi problemi a tutto il sistema finanziario ed economico in caso di crisi o fallimento.

Come si evince dalla Tabella 1, il credito è di gran lunga la principale attività per le banche etiche (76,11% del totale nel 2018), mentre rappresenta poco più di un terzo delle attività (39,80% nel 2018) per il sistema bancario europeo. Le banche etiche si confermano, quindi, più votate all'attività bancaria classica (raccolta di risparmi e concessione di crediti) rispetto al settore bancario europeo, che appare invece concentrato su altri tipi di attività quali investimenti in titoli, servizi finanziari, partecipazioni in imprese e altro.

Dal 2008 al 2018 il peso dell'attività creditizia sul totale è cresciuto in valore assoluto per entrambi gli aggregati, anche se in misura minore per il gruppo delle "banche europee" (+2,68 punti percentuali) rispetto alle banche etiche e sostenibili (+3,83 punti percentuali). Se si guardano gli incrementi relativi dell'ultimo decennio, che sono anche gli anni che seguono la crisi finanziaria del 2007-2008, le dinamiche si invertono, segnando un 7,2% per le prime e un 5,3% per le seconde. Visto che il credito può essere considerato, con qualche approssimazione, un'attività di finanziamento dell'economia reale (in assenza di dati più precisi nei bilanci delle banche), è possibile concludere che le banche etiche e sostenibili operano decisamente a sostegno dell'economia reale (produzione di beni e servizi tangibili) mentre il resto del sistema bancario europeo, in media, è più orientato all'economia finanziaria (investimenti in borsa, vendita di titoli, ecc.).

Proseguendo nell'analisi, la differenza tra i due gruppi di banche è confermata anche dal rapporto percentuale tra depositi e totale del Passivo. La Tab. 2 evidenzia come le banche etiche e sostenibili raccolgano denaro (che poi, principalmente, prestano in forma di crediti) soprattutto tramite i depositi dei clienti (71,31% del totale) mentre, in media, le banche europee raccolgono liquidità (da prestare o investire) soprattutto da altri canali, come per esempio con emissione di obbligazioni o depositi di altre banche. Solo il 40,96% del passivo delle banche europee è costituito da depositi: una percentuale che è cresciuta di circa otto punti dal 2008 in corrispondenza, però, di una crescita altrettanto sostenuta per le banche etiche e sostenibili. In termini di variazione relativa, l'aggregato delle banche etiche/sostenibili europee è cresciuto dell'11,93%, a fronte del +31,03% del sistema nel suo complesso.

Tabella 2. *Peso percentuale dei depositi sul totale delle passività rispetto alla tipologia di banche (etiche vs complesso)*

Tipologia di banche rispetto alla caratteristica della eticità	Depositi/Totale passivo		
	2008	2013	2018
Banche etiche/sostenibili europee	63,71	76,95	71,31
Banche europee	31,26	36,57	40,96

Fonte: *Fondazione Banca Etica*

Le banche etiche e sostenibili possono vantare una solida posizione patrimoniale, misurata come rapporto tra il patrimonio netto e il totale del passivo, mantenuta costante dal 2008 al 2018 e attestata intorno al 10%. Il sistema bancario europeo nel suo complesso, invece, è partito da una posizione relativamente più debole nel 2008 (5,55%) per poi colmare progressivamente lo scarto con le banche etiche, (8,18% nel 2018), con una variazione relativa di poco inferiore al 50%. Questo perché le banche europee “tradizionali” sono state spinte dal legislatore a incrementare il proprio patrimonio per far fronte a possibili future crisi dopo quella del 2007-2008, mentre le banche etiche hanno sempre avuto un patrimonio netto relativamente elevato rispetto al totale del passivo, anche nel periodo che precede il 2007-2008 (Tabella 3).

Tabella 3. *Peso percentuale del patrimonio netto sul totale delle passività rispetto alla tipologia di banche (etiche vs complesso)*

Tipologia di banche rispetto alla caratteristica della eticità	Patrimonio netto/Passivo		
	2008	2013	2018
Banche etiche/sostenibili europee	10,05	9,90	10,54
Banche europee	5,55	7,89	8,18

Fonte: Fondazione Banca Etica

Passando all’analisi reddituale, essa viene condotta attraverso i tradizionali indici di bilancio ROA (Return on Assets) e ROE (Return on equity), calcolati per le banche etiche europee e per il sistema bancario europeo. Come è noto, il ROA è ottenuto come rapporto tra l’utile netto e il totale dell’attivo ed è una misura della redditività delle attività. La tabella 4 illustra i dati di media e deviazione standard dei rispettivi ROA per i due gruppi di banche nell’ultimo quinquennio disponibile e nel periodo 2008-2018.

Tabella 4. *Principali indici di performance (media e deviazione standard) per il ROA-Return On Assets a 5 e 10 anni*

Tipologia di banche rispetto alla caratteristica della eticità	ROA - RETURN ON ASSETS	
	5 anni (2013-2018)	
	Media	Deviazione standard
Banche etiche/sostenibili europee	0,40%	0,05%
Banche europee	0,26%	0,13%
	10 anni (2008-2018)	
	Media	Deviazione standard
	Banche etiche/sostenibili europee	0,40%
Banche europee	0,13%	0,20%

Fonte: Fondazione Banca Etica

Nel periodo 2008-2018, il ROA medio delle banche etiche e sostenibili si è mantenuto sempre su livelli superiori (in media 0,40%) rispetto a quello del sistema bancario europeo, con una volatilità, misurata mediante la deviazione standard calcolata rispetto al valore medio di ogni anno, decisamente bassa (0,05% nel periodo 2013-2018 e 0,40% sull'arco temporale più lungo). In entrambi i periodi analizzati, il sistema bancario europeo ha registrato una redditività media inferiore rispetto alle banche etiche e sostenibili (0,13% vs 0,40% sui dieci anni), con una volatilità decisamente maggiore.

Il ROE (Return on Equity) è il rapporto tra l'utile netto e il patrimonio netto ed è una misura del rendimento contabile di un'impresa.

Come illustrato in Tabella 5, si ripropone la situazione già registrata per il ROA: il rendimento medio delle banche etiche è stato più alto rispetto a quello del sistema bancario europeo nel periodo 2008-2018 (3,57% vs 1,79%) associato ad una volatilità, e quindi un livello di rischio, molto inferiore (0,41% vs 3,28%). Tuttavia, negli ultimi cinque anni (2013- 2018), le banche europee sono tornate a crescere in termini di rendimenti, superando, anche se di poco (3,86% vs 3,61%), le banche etiche e sostenibili, a costo, però, di un rischio più che triplicato (1,72% vs 0,48%).

Tabella 5. *Principali indici di performance (media e deviazione standard) per il ROE-Return On Equity a 5 e 10 anni*

Tipologia di banche rispetto alla caratteristica della eticità	ROE-RETURN ON EQUITY	
	5 anni (2013-2018)	
	<i>Media</i>	<i>Deviazione standard</i>
Banche etiche/sostenibili europee	3,61%	0,48%
Banche europee	3,86%	1,72%
	10 anni (2008-2018)	
	<i>Media</i>	<i>Deviazione standard</i>
	Banche etiche/sostenibili europee	3,57%
Banche europee	1,79%	3,28%

Fonte: Fondazione Banca Etica

Per completezza di analisi è interessante osservare i dati percentuali di crescita delle principali voci di bilancio (attivi, prestiti, depositi e patrimonio netto) per entrambi i gruppi di banche. Tale variabile è stata calcolata come tasso annuo di crescita composto, più comunemente noto come CAGR dall'acronimo anglosassone Compounded Average Growth Rate, che rappresenta la crescita percentuale media di una grandezza nel periodo di tempo preso in analisi.

Come evidenziato in Tabella 6, negli ultimi dieci anni, le banche etiche e sostenibili sono cresciute molto di più rispetto al sistema bancario europeo. In particolare, esse hanno raggiunto dimensioni maggiori, come attestato dalla crescita annua degli attivi del 9,91% negli ultimi dieci anni, mentre le banche europee si sono ridimensionate (-0,31% in media all'anno negli ultimi dieci anni); hanno concesso un volume maggiore di prestiti (+10,55% all'anno dal 2008 al 2018 contro il +0,39% delle banche europee) e raccolto più denaro tramite depositi (+11,17% vs +2,43%). Infine, hanno aumentato in modo considerevole il proprio patrimonio netto, +10,40% in media all'anno negli ultimi dieci anni, contro il +3,65% delle banche europee.

Tabella 6. *Principali indici della crescita del collettivo delle banche a 5 e 10 anni.*

Tipologia di banche rispetto alla caratteristica della eticità	Indici di crescita	
	5 anni (2013-2018)	10 anni (2008-2018)
	Totale attivo	
Banche etiche/sostenibili europee	7,94%	9,91%
Banche europee	0,31%	-0,31%
	Prestiti	
Banche etiche/sostenibili europee	8,69%	10,55%
Banche europee	0,88%	0,39%
	Depositi	
Banche etiche/sostenibili europee	6,31%	11,17%
Banche europee	2,61%	2,43%
	Patrimonio netto	
Banche etiche/sostenibili europee	9,30%	10,40%
Banche europee	1,04%	3,65%

Fonte: *Fondazione Banca Etica*

E' utile osservare che, se dopo l'ultima grande crisi finanziaria del 2007-2008, il sistema bancario europeo ha vissuto, in generale, un lungo periodo di stagnazione o di crescita molto debole, con un trend che si è invertito solo negli ultimi cinque anni, le banche etiche e sostenibili hanno continuato a crescere in modo rilevante, probabilmente sostenute anche dal desiderio di molti risparmiatori di trovare un'alternativa alle banche tradizionali, molte delle quali sono state tra le principali responsabili, ma per certi versi anche vittime, della crisi.

4. Considerazioni conclusive

Sin dalla loro fondazione, le banche etiche sono risultate molto più orientate a offrire servizi all'economia reale rispetto alla media del sistema bancario, vantando una maggior solidità dal punto di vista patrimoniale e una maggiore redditività per entrambi gli indicatori utilizzati (ROA e ROE) a fronte di una minore volatilità, e quindi di minori rischi.

Tuttavia, negli ultimi cinque anni oggetto di analisi (2013-2018), le banche europee sono tornate a crescere in termini di rendimenti, superando, anche se di poco, le banche etiche e sostenibili. Mentre, in media, il sistema bancario europeo sembra essere uscito dalla crisi per imboccare un sentiero di crescita progressiva della redditività, le banche etiche e sostenibili continuano ad avere rendimenti pressoché costanti, con una leggera flessione nel 2018, dovuta in gran parte al crollo dell'utile netto di una sola banca, il *Crédit Coopératif*.

Le banche etiche hanno, inoltre, registrato una crescita considerevole di tutte le grandezze misurate dall'analisi negli ultimi dieci anni, mentre il sistema bancario europeo è cresciuto in misura molto minore. Le future analisi dovranno necessariamente mettere a confronto due periodi temporali differenti, definibili "ante e post-Covid19".

A tal proposito, occorre evidenziare come la pandemia da Coronavirus abbia esponenzialmente focalizzato l'attenzione sulle tematiche ambientali, sociali e di governance. Infatti, gli investitori riconoscono che questi c.d. fattori non finanziari, tra i quali aspetti di salute e di sicurezza delle aziende, il rapporto con i dipendenti, la capacità dei Consigli di Amministrazione di gestire questi temi, ecc., saranno estremamente importanti nel determinare quali banche potranno emergere da questa crisi in modo soddisfacente, laddove quelle con elevati standard dal punto di vista dell'Environment, Social and Governance (ESG) hanno maggiori possibilità di trattenere forza lavoro motivata, clienti fedeli e linee di produzione resilienti.

I primi segnali in tal senso sembrano incoraggianti: come la recente indagine HSBC 2020 "Sustainable Financing and Investing", condotta su un campione di circa 9.000 aziende ha evidenziato, i protagonisti del mercato europeo dei capitali attribuiscono una importanza ben più elevata alla finanza sostenibile rispetto al periodo precedente la pandemia di COVID-19. L'indagine riporta che, secondo il 36% degli emittenti europei, la pandemia ha portato in primo piano le questioni legate alla sostenibilità e, inoltre, il 77% degli intervistati ha dichiarato che la

pandemia li ha indotti ad incrementare il loro impegno verso i temi ESG e verso la componente sociale degli stessi.

Attualmente, gli investitori europei sembrano mostrare la maggior sensibilità verso le questioni ambientali e sociali rispetto ai competitors mondiali, e risultano più predisposti a considerare il ruolo dei fattori ESG nelle loro decisioni di investimento (39% delle aziende europee vs il 31% delle aziende che operano a livello mondiale). Tra gli ambiti di maggior interesse degli investitori europei rientrano: le opportunità relative alle infrastrutture sostenibili, le energie rinnovabili, le smart cities. La ricerca ha, infine, evidenziato che le azioni delle aziende con rating ESG più elevati hanno sovraperformato la media globale del 4,7% a partire da metà dicembre 2019. Tale divario nella performance diviene più ampio per i titoli legati al clima, che hanno registrato una performance superiore del 13% rispetto alla media globale.

La finanza sostenibile è, ormai, uno dei mercati dei capitali che si sta sviluppando più velocemente e gli accordi ESG vengono considerati tra gli asset tradizionali. Secondo gli analisti, il successo delle aziende nel lungo periodo sarà determinato dalla loro capacità di effettuare investimenti ESG competitivi rispetto agli investimenti tradizionali in termini di rapporto rischio-rendimento. L'Europa, come anticipato, è uno dei più forti sostenitori internazionali della necessità di combattere il cambiamento climatico ed è leader mondiale della finanza sostenibile; conseguentemente, da una parte, gli investitori europei sono chiamati a svolgere un ruolo chiave e, dall'altra, i cittadini devono continuare ad investire ispirati dai principi della sostenibilità ambientale.

Bibliografia essenziale

- Angelini, A., Farioli, F., Mattili, F. (2015). Le due crisi. Crisi del capitalismo e crisi ambientale: una soluzione sostenibile? *Culture della sostenibilità*, 16: 95-114,
- Antonucci, A. (2007). La responsabilità sociale d'impresa. *La Nuova Giurisprudenza Civile Commentata*, 2: 119-129.
- Barnett, M., Salomon R. (2006). Beyond Dichotomy: The Curvilinear Relationship Between Social Responsibility and Financial Performance. *Strategic Management Journal*, 27: 1101-1122.
- Beal, D., Goyen, M., Phillips, P. (2005). Why do we invest ethically? *Journal of Investing*, 14(3): 66-77.

- Cavallito, M. Isonio, E., Meggiolaro, M. (2020). *Terzo rapporto La finanza etica e sostenibile in Europa*, Fondazione Finanza Etica, Collana Studi e Ricerche.
- Commissione Europea (2001). *Libro verde – Promuovere un quadro europeo per la responsabilità sociale delle imprese*, COM 2001/398, 18 luglio.
- Chava, S. (2014). Environmental externalities and cost of capital. *Management Science*, 60(9): 2223-2247.
- Dal Maso, D., Fiorentini, G. (2013). *Creare valore a lungo termine: Conoscere, promuovere e gestire l'investimento sostenibile e responsabile*. Milano, EGEA
- Desmadryl, X., (2007), *SRI & ESG inclusion: does it pay after all?*, ANBID/UNEP FI roundtable in São Paolo, Brazil, for HSBC, March 2007.
- Eccler, G., Ioannou I., Serafiem, G. (2013). *The impact of corporate sustainability on organizational processes and performance*, Harvard Business School working papers.
- Fulton, M., Kahn, B. and Sharples, C. (2012). Sustainable investing: Establishing long-term value and performance. Available at SSRN: <https://ssrn.com/abstract=2222740> .
- Galbraith, J.K. (2000). *Breve storia dell'euforia finanziaria*, Rizzoli.
- Gallino, L. (2005). *L'impresa irresponsabile*, Einaudi.
- Global Reporting Initiative (2016), *Linee guida per il reporting di sostenibilità*.
- Kreander N. Gray, R. H. Power, D.M., Sinclair, C.D. (2005) Evaluating Performance of Ethical and Non-ethical Funds: A Matched Pair Analysis. *Journal of Business Finance & Accounting*, Vol. 32, No. 7-8: 1465-1493.
- Perrini, F., Tencati, A. (2011). *Corporate social responsibility: Un nuovo approccio strategico alla gestione d'impresa*. EGEA Milano.
- Regalli, M., Soana, M. G, Tagliavini, G. (2005). I fondi etici: caratteristiche, spazi di mercato, ritorni finanziari. *Etica e Finanza*, FrancoAngeli, Milano: 177-200.
- Stiz, G., Seme, C. I. (1999). *Guida alla finanza etica. Come investire i propri risparmi in modo socialmente responsabile*, Editrice Missionaria Italiana 2000, 16-17.
- Tomasi F., Russo A. (2012). *Il rating etico: un'analisi empirica del modello standard ethics*, Pubblicazioni Standard Ethics Research.
- Valor, C., De la Cuesta, M. (2007). Códigos éticos: análisis de la eficacia de su implantación entre las empresas españolas cotizadas. *INNOVAR. Revista de Ciencias Administrativas y Sociale*, Universidad Nacional de Colombia Editor, vol.17, 30: 19-30..

Sitografia

www.bancaetica.it

www.eurosif.org

www.finanzasostenibile.it

www.borsaitaliana.it

www.climatebonds.net

www.economicircolare.com

www.gbm.hsbc.com/solutions/sustainable-financing

Analisi delle strutture evolutive imprenditoriali a livello regionale e comunale: una applicazione su Puglia e Basilicata

Agata Maria Madia Carucci¹, Giovanni Vannella^{2*}.

¹Istat, ²Università degli Studi di Bari.

Riassunto: Nel presente lavoro viene ad essere illustrata una analisi interpretativa delle evoluzioni competitive imprenditoriali condotta a livello territoriale, con dettaglio comunale, tramite una applicazione dell'analisi Shift-share utilizzata sui dati provenienti dall'archivio ASIA Unità Locali e pubblicati nella banca dati I.stat.

Keywords: Statistiche strutturali a livello comunale e regionale, analisi competitività, analisi Shift-share.

1. Introduzione

L'analisi delle competitività territoriali e, più in generale, delle variazioni spazio-temporali delle variabili economiche, è da sempre oggetto di studi ed approfondimenti mirati ad individuare come e perché tali differenziazioni si manifestino e si diversifichino, unitamente alle conseguenze che le stesse diversità generino. Ovviamente tali analisi risentono fortemente delle scelte localizzative e delle capacità di evolversi delle imprese.

Le motivazioni dei diversi dinamismi territoriali possono ascrivere ad una serie di fattori tra cui si veda, ad esempio, la differente dotazione infrastrutturale e di capitale umano del territorio, la presenza di imprese con cui interagire in modo profittevole, la maggiore o minore efficienza della PA e di politiche attrattive, il minore o maggiore spirito imprenditoriale dell'area, etc.

Prima di capire quali possano essere le cause, occorre però verificare quali differenze esistano tra i territori e, al fine di tale verifica, un aiuto importante

* Autore corrispondente: giovanni.vannella@uniba.it. Sebbene il presente articolo sia frutto del lavoro congiunto dei due autori, a A.M.M. Carucci sono attribuiti i par. 2 e 3, a G. Vannella i par. 1, 4, 5.

proviene dalla statistica ufficiale che permette di ottenere informazioni puntuali e celeri sulle strutture evolutive delle realtà imprenditoriali dei territori.

Negli ultimi tempi gli istituti centrali di statistica riescono ad assolvere a tale compito con crescente capacità di coniugare celerità dell'offerta informativa unitamente ad una sempre maggiore puntualità territoriale della stessa, il tutto grazie al crescente maggiore utilizzo delle fonti amministrative integrate (Seri *et al.*, 2016) quali il frame-sbs (Boselli *et al.*, 2016) che ha portato alla definizione di un nuovo approccio di stima delle variabili economiche su base territoriale (Faramondi *et al.*, 2018).

Già in passato, gli autori del presente lavoro si sono occupati di aspetti particolari relativi a tale processo di integrazione tra fonti (Carucci *et al.*, 2016) verificando le potenzialità del sistema informativo "a misura di comune" (Bianchino *et al.*, 2018) e del frame-sbs territoriale (Carucci *et al.*, 2019).

In tali lavori erano emerse le notevoli potenzialità dei nuovi sistemi informativi nonché l'opportunità di condurre, su tali fonti, analisi del dinamismo locale attraverso semplici applicazioni ed adattamenti dell'analisi Shift-share.

Si era inoltre evidenziato come potesse essere interessante verificare la fattibilità di tale approccio di analisi con riferimento ad ambiti territoriali molto dettagliati (province e comuni).

Alla luce di ciò si è voluta verificare la potenzialità dell'analogo approccio di analisi Shift-share sui dati provenienti dall'archivio ASIA unità locali, diffusi tramite il portale I.stat, le cui risultanze (in parte discusse dagli stessi autori al recente congresso AISRE 2020) vengono illustrate nel presente lavoro.

2. Le statistiche sulle imprese presenti nella banca dati I.Stat

La banca dati I.Stat rappresenta il data warehouse dell'Istat nonché la più completa forma di offerta statistica dell'Ente.

I dati sono liberamente accessibili e vengono rappresentati con una suddivisione che prevede una organizzazione per temi, presentati in tavole multidimensionali e corredati di un ampio apparato di metadati.

Per quanto concerne le imprese, la banca dati fornisce una rilevante serie di informazioni sui seguenti sei macro temi: struttura, competitività, innovazione nelle imprese con almeno 10 addetti, ICT nelle imprese con almeno 10 addetti, ricerca e sviluppo, clima di fiducia delle imprese.

I dati sulla struttura sono suddivisi nelle seguenti sotto-aree:

- Imprese e addetti
- Imprese - occupati
- Imprese - occupati per titolo di studio
- Imprese - lavoratori dipendenti
- Imprese - lavoratori dipendenti beneficiari di sgravi contributivi
- Imprese - lavoratori indipendenti
- Imprese - lavoratori esterni
- Imprese - lavoratori esterni per compensi percepiti e periodo di attività
- Imprese - lavoratori temporanei
- Unità locali e addetti
- Unità locali - lavoratori dipendenti.

Essendo l'analisi condotta incentrata su una valutazione a livello territoriale, si è ritenuto preferibile lavorare sulle unità locali (di seguito UL), piuttosto che sulle imprese, in quanto più idonee, per ovvi motivi, a rappresentare le peculiarità locali. La struttura informativa di I.Stat per quanto concerne le UL, disponibile per gli anni dal 2012 in poi, è così organizzata:

- Unità locali e addetti
 - Principali dati
 - Settori economici (Ateco 5 cifre) – tipo e ampiezza demografica dei comuni e province
 - Settori economici (Ateco 3 cifre) – comuni
 - Classe di addetti, settori economici (Ateco 5 cifre) – Italia
 - Classe di addetti, settori economici (Ateco 4 cifre) – ripartizioni
 - Classe di addetti, settori economici (Ateco 3 cifre) – province
 - Classe di addetti, settori economici (Ateco 2 cifre) – comuni
 - Dati per sistema locale del lavoro
 - Classe di addetti, settori economici (Ateco 1 cifra) – sistemi locali del lavoro 2011
 - Classe di addetti, settori economici (Ateco 1 cifra) – sistemi locali del lavoro 2001
- Unità locali - lavoratori dipendenti
 - Classe di addetti, sesso, età – province
 - Classe di addetti, paese di nascita – province
 - Classe di addetti, qualifica professionale – province
 - Settori economici (Ateco 3 cifre), sesso – province

- Settori economici (Ateco 3 cifre), età – province
- Settori economici (Ateco 3 cifre), paese di nascita – province
- Settori economici (Ateco 3 cifre), qualifica professionale – province.

I dati sulle imprese derivano dal Registro statistico delle imprese attive ASIA, mentre quelli delle UL dal Registro statistico delle UL (ASIA - UL). I due registri presentano oramai una struttura consolidata ed ampiamente testata traendo entrambi origine dal Regolamento del Consiglio Europeo n. 2816/93 relativo al coordinamento comunitario dello sviluppo dei registri d'impresa utilizzati a fini statistici, poi abrogato e sostituito dal Regolamento CE n. 177/2008.

Presentano lo stesso campo di osservazione costituito dalle unità economiche relative alle attività industriali, commerciali e dei servizi alle imprese ed alle famiglie. Non sono però incluse alcune tipologie di attività economiche, ovvero quelle relative alle sezioni Ateco A (agricoltura, silvicoltura e pesca), O (amministrazione pubblica e difesa, assicurazione sociale obbligatoria), T (produzione di beni e servizi indifferenziati per uso proprio da parte di famiglie e convivenze), U (organizzazioni ed organismi extraterritoriali), e quelle relative alla sezione S, divisione 94 (attività di organizzazioni associative) nonché tutte le attività classificate come istituzioni pubbliche e istituzioni private non profit.

Per quanto concerne in dettaglio il Registro ASIA-UL, da un punto di vista definitorio l'unità locale corrisponde ad un'impresa o a una parte di un'impresa situata in una località topograficamente identificata in cui una o più persone svolgono attività economiche (non necessariamente a tempo pieno) per conto di una stessa impresa (regolamento del Consiglio Europeo N. 696 del 15 marzo 1993). Ne discende che l'impresa plurilocalizzata svolge le proprie attività in più luoghi, ciascuno dei quali costituisce un'unità locale.

Le variabili specifiche delle UL comprese nel Registro, oltre alle variabili identificative dell'impresa e definite nel Registro ASIA-Imprese, sono:

- indirizzo dell'UL, che permette l'esatta individuazione dell'unità locale sul territorio;
- Attività economica dell'UL, secondo la classificazione Ateco 2007;
- Addetti dell'unità locale.

Il Registro delle UL prevede un periodico annuale percorso di aggiornamento attraverso l'integrazione di fonti amministrative e statistiche.

3. Il contesto delle strutture imprenditoriali lucane e pugliesi

Prima di procedere con l'applicazione dell'analisi Shift-share, si è ritenuto opportuno analizzare la distribuzione e l'andamento, per le regioni Puglia e Basilicata e per il periodo 2012-2018, degli aggregati "addetti" ed "unità locali", poi utilizzati per tale analisi, elaborando i dati I.stat.

Una particolare attenzione è stata rivolta alla composizione per attività economica delle unità produttive presenti sul territorio confrontandole, laddove necessario, con le corrispondenti risultanze nazionali e della macro ripartizione territoriale.

In Italia, a fronte di una variazione negativa delle UL dell'1,4% nell'intero periodo in esame, gli addetti sono aumentati, del 3,4%. La ripartizione Sud e le due regioni, Puglia e Basilicata, presentano un leggero incremento delle UL, che si riflette su un aumento più significativo degli addetti. È la Basilicata a registrare l'incremento, in termini di addetti, più alto, pari ad oltre il 10% (Tab. 1).

Tabella 1. *Addetti in Italia, nella ripartizione Sud e nelle regioni Puglia e Basilicata. Variazione percentuale del 2018 su 2012.*

Aree territoriali	Variazione 2018/2012	
	Addetti	Unità locali
Italia	3,4	-1,4
Sud	4,8	0,2
Puglia	4,5	0,1
Basilicata	10,2	0,5

Elaborazione su dati I.stat

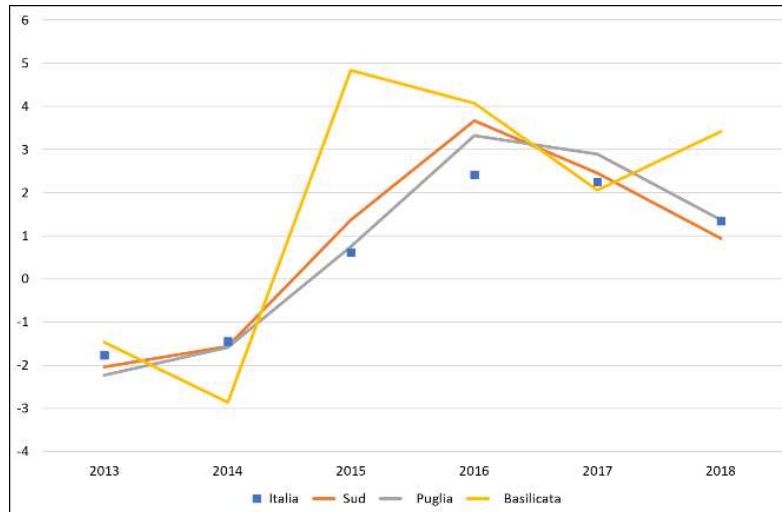
La variazione degli addetti e delle UL nell'intero periodo è la sintesi dell'andamento dei due aggregati in ciascun anno rispetto all'anno precedente, andamento che, come facilmente immaginabile, non si è rilevato essere costante nel periodo considerato.

Infatti, come si vede in Fig. 1, in Italia sino al 2016, per quanto concerne gli addetti, si registra una crescita cui si contrappone una decrescita negli anni successivi.

Lo stesso aggregato si ripresenta in modo sostanzialmente omogeneo con riferimento al Sud Italia ed alla Puglia, mentre la Basilicata presenta un andamento leggermente diverso con un valore minimo del 2014 per poi crescere più della media nazionale e ripartizionale nell'anno successivo; nell'ultimo anno si registra una ulteriore variazione positiva, con una crescita degli addetti del 3,4% nel 2018 rispetto al 2017.

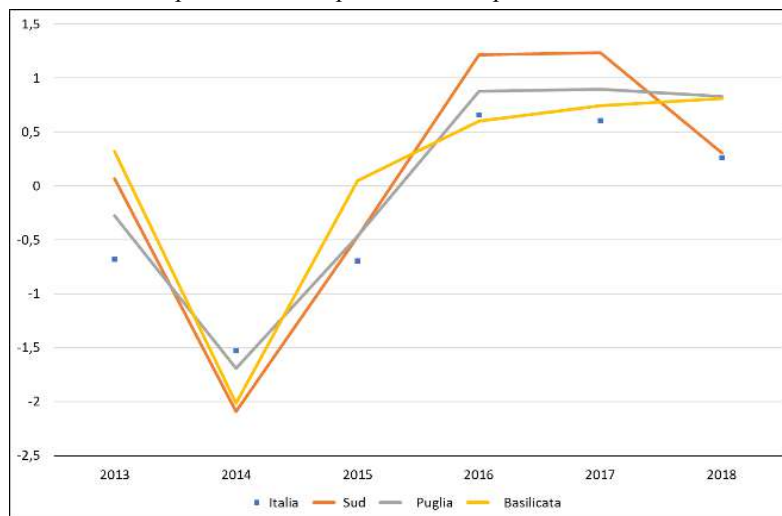
In termini di UL la variazione nelle regioni non ha difformità significative da quella ripartizionale e nazionale (Fig. 2).

Figura 1. Addetti in Italia, nella ripartizione Sud e nelle regioni Puglia e Basilicata. Variazione percentuale rispetto all'anno precedente. Anni 2013-2018.



Elaborazione su dati I.stat

Figura 2. Unità locali in Italia, nella ripartizione Sud e nelle regioni Puglia e Basilicata. Variazione percentuale rispetto all'anno precedente. Anni 2013-2018.



Elaborazione su dati I.stat

Passando all'analisi disaggregata a livello provinciale (Tab. 2), in Basilicata, Matera registra nell'intero periodo un incremento di circa il 2% delle UL, che si traduce in oltre il 12% di addetti in più. Nella provincia di Potenza all'incremento degli addetti corrisponde una variazione di UL pressoché nulla.

In Puglia l'incremento maggiore, sia in termini di addetti che di UL, si ha nelle province di Barletta-Andria-Trani e Bari. Le province con le performance peggiori sono Foggia e Taranto.

Tabella 2. *Addetti nelle province di Puglia e Basilicata. Variazione percentuale del 2018 su 2012.*

Aree territoriali	Variazione 2018/2012	
	Addetti	Unità locali
Basilicata	10,2	0,5
Potenza	9,1	-0,3
Matera	12,6	2,0
Puglia	4,5	0,1
Foggia	2,9	-1,3
Bari	6,8	0,7
Taranto	0,1	-0,4
Brindisi	3,8	0,4
Lecce	3,3	-0,1
Barletta-Andria-Trani	7,6	1,4
Basilicata	10,2	0,5
Potenza	9,1	-0,3
Matera	12,6	2,0

Elaborazione su dati I.stat

Passando alla disaggregazione della variazione per anno degli addetti (Tab. 3) e delle UL (Tab. 4) si osserva come, con riferimento alla Basilicata, insista un andamento provinciale degli addetti molto diverso nelle due province: in particolar modo le differenze più significative si hanno tra il 2015 e il 2018. Nel biennio 2015-2016 è la provincia di Potenza a crescere maggiormente rispetto all'anno precedente, mentre nel successivo biennio la provincia di Matera cresce di oltre 5 punti percentuali a fronte di una crescita molto più contenuta osservata per la provincia di Potenza. In termini di UL, invece, a partire dal 2015, è la provincia di Matera a registrare livelli di crescita più alti.

In Puglia le variazioni nel tempo dei due aggregati appaiono meno facilmente leggibili. Con riferimento infatti al confronto tra performance provinciali e regionali, Foggia presenta una pressoché costante sotto performance per entrambi gli aggregati, Barletta-Andria-Trani presenta quasi sempre un andamento prevalentemente più elevato di quello regionale per entrambi gli aggregati, le province di Brindisi, Lecce e quella di Bari tendono prevalentemente a ricalcare l'andamento medio regionale sebbene Bari su livelli leggermente più alti dal 2015, Taranto tende infine a registrare una prevalenza di sotto performance.

Tabella 3. *Addetti in Italia, nella ripartizione Sud e nelle regioni Puglia e Basilicata. Variazione percentuale rispetto all'anno precedente. Anni 2013-2018.*

Aree territoriali	2013	2014	2015	2016	2017	2018
Basilicata	-1,5	-2,9	4,8	4,1	2,1	3,4
Potenza	-0,4	-3,7	6,2	4,4	0,7	1,9
Matera	-3,6	-1,1	2,0	3,4	5,1	6,6
Puglia	-2,2	-1,6	0,8	3,3	2,9	1,4
Foggia	-2,5	-1,4	0,0	4,9	1,8	0,2
Bari	-2,3	-1,4	1,3	4,2	3,2	1,7
Taranto	-1,8	-3,5	2,1	0,9	1,2	1,4
Brindisi	-3,1	-0,8	1,0	2,6	4,3	-0,1
Lecce	-3,0	-1,4	-0,3	2,3	3,1	2,7
Barletta-Andria-Trani	0,0	-0,9	-0,3	4,1	4,1	0,4

Elaborazione su dati I.stat

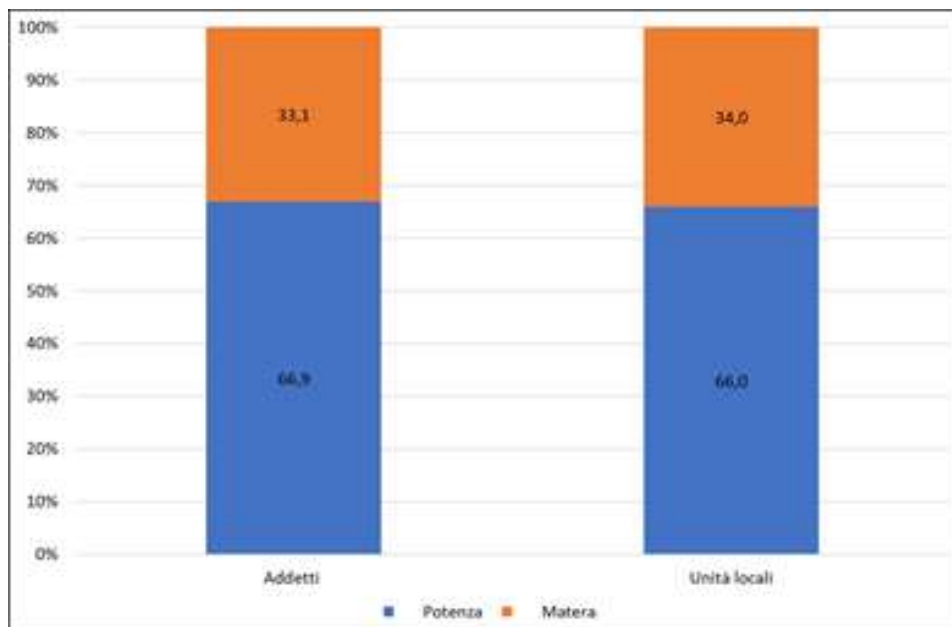
Tabella 4. *Unità locali in Italia, nella ripartizione Sud e nelle regioni Puglia e Basilicata. Variazione percentuale rispetto all'anno precedente. Anni 2013-2018.*

Aree territoriali	2013	2014	2015	2016	2017	2018
Basilicata	0,3	-2,0	0,0	0,6	0,7	0,8
Potenza	0,6	-2,2	-0,2	0,4	0,6	0,6
Matera	-0,3	-1,6	0,5	1,1	1,0	1,3
Puglia	-0,3	-1,7	-0,5	0,9	0,9	0,8
Foggia	-0,2	-1,9	-1,1	1,0	0,4	0,5
Bari	-1,1	-1,3	0,2	0,8	0,7	1,5
Taranto	0,0	-2,9	-0,2	1,4	1,0	0,3
Brindisi	0,0	-1,8	-0,6	0,9	1,2	0,7
Lecce	-0,7	-1,7	-0,7	0,9	1,1	1,1
Barletta-Andria-Trani	3,0	-1,3	-1,4	0,3	1,5	-0,6

Elaborazione su dati I.stat

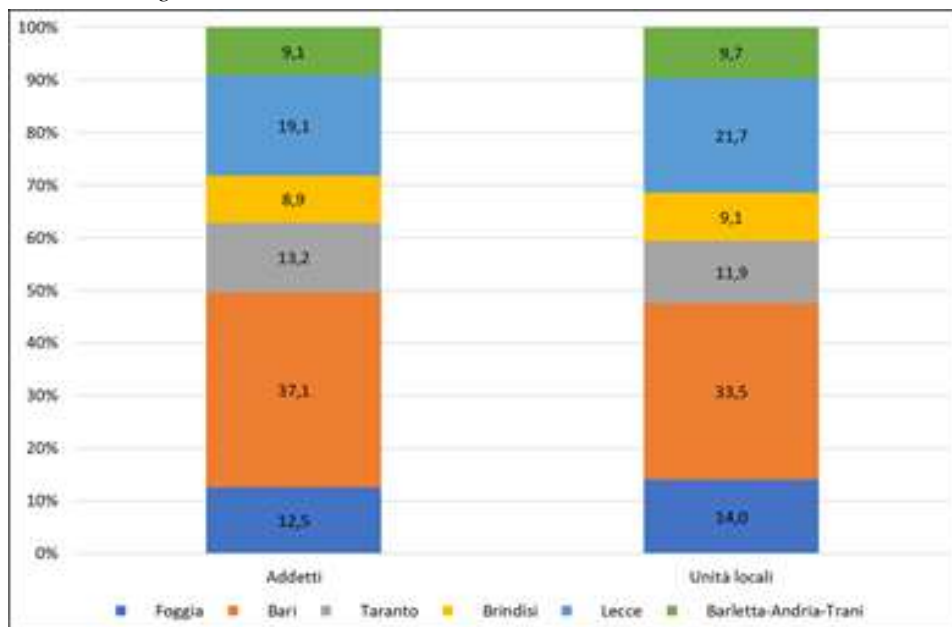
All'interno della regione, ciascuna provincia contribuisce in maniera diversa alla variazione degli addetti e delle UL e tale contributo è strettamente legato al peso relativo di ciascuna provincia nella regione stessa e alla composizione per attività economica degli addetti e delle UL di ciascuna provincia. Per la Basilicata, la provincia di Potenza, per esempio, pesa nel 2018 per oltre il 60% sul totale addetti ed UL regionale (Fig. 3). Tale peso è costante nell'intero periodo in esame e spiega l'andamento simile della provincia di Potenza e della regione lucana. Sia in termini di addetti che di UL, in Puglia, è la provincia di Bari a pesare di più (oltre il 30%), segue la provincia di Lecce (circa il 20%), mentre il peso minore per addetti e UL è della provincia di Brindisi (circa il 9%) (Fig. 4).

Figura 3. *Composizione percentuale per provincia degli addetti e delle unità locali in Basilicata. Anno 2018.*



Elaborazione su dati I.stat

Figura 4. *Composizione percentuale per provincia degli addetti e delle unità locali in Puglia. Anno 2018.*



Elaborazione su dati I.stat

Tabella 5. *Composizione percentuale degli addetti e delle unità locali per attività economica in Italia e nella ripartizione Sud. Anno 2018.*

Attività economica	Italia		Sud	
	Addetti	Unità locali	Addetti	Unità locali
Industria	23,3	9,5	19,2	8,7
Costruzioni	7,5	10,6	8,8	9,7
Servizi, commercio, trasporto, alloggio e ristorazione	35,4	35,8	40,9	42,9
Servizi di informazione e comunicazione	3,3	2,5	2,2	1,8
Attività finanziarie e assicurative	3,1	2,8	2,3	2,4
Attività immobiliari	1,7	5,0	0,9	2,4
Attività professionali, scientifiche, tecniche ed amministrative	15,6	20,1	14,4	18,9
Istruzione, salute e servizi sociali	6,0	7,5	6,8	7,0
Altri servizi	3,9	6,2	4,4	6,1
Totale	100,0	100,0	100,0	100,0

Elaborazione su dati I.stat

Tabella 6. *Composizione percentuale degli addetti e delle unità locali per attività economica in Basilicata e province. Anno 2018.*

Attività economica	Basilicata		Potenza		Matera	
	Addetti	Unità locali	Addetti	Unità locali	Addetti	Unità locali
Industria	26,6	9,6	28,1	9,9	23,5	9,1
Costruzioni	10,2	11,1	10,5	11,5	9,6	10,2
Servizi, commercio, trasporto, alloggio e ristorazione	35,1	41,8	33,0	40,9	39,4	43,4
Servizi di informazione e comunicazione	2,2	1,9	2,5	1,9	1,8	1,9
Attività finanziarie e assicurative	2,0	2,4	2,0	2,5	2,1	2,1
Attività immobiliari	0,5	1,5	0,5	1,7	0,5	1,3
Attività professionali, scientifiche, tecniche ed amministrative	13,9	19,2	14,5	19,3	12,8	19,1
Istruzione, salute e servizi sociali	6,0	6,5	5,8	6,4	6,3	6,6
Altri servizi	3,4	6,0	3,1	5,9	4,1	6,3
Totale	100,0	100,0	100,0	100,0	100,0	100,0

Elaborazione su dati I.stat

Passando ora al confronto delle suddivisioni delle attività imprenditoriali per attività economica (Tab. 5), in media in Italia nel 2018 i più importanti macrosettori di attività con riferimento agli addetti sono nell'ordine "servizi", "industria" e "attività professionali", mentre con riferimento alle UL sono "servizi", "attività professionali" e "costruzioni".

Tale struttura si ripresenta anche con riferimento alla ripartizione Sud pur se con un maggior peso dei "servizi" ed un minore peso dell'"industria".

In Basilicata è più alto, rispetto alla media nazionale, il peso dei macrosettori "industria" e "costruzioni", mentre in Puglia quello dei "servizi".

Con riferimento ai contesti provinciali (Tab. 6), in Basilicata, a Potenza si rileva una presenza maggiore del macrosettore "industria" (28% di addetti e 10% di UL) mentre a Matera del macrosettore "servizi" (39% di addetti e 43% di UL).

In Puglia spicca invece la forte vocazione industriale con riferimento agli addetti di Barletta-Andria-Trani e Taranto (Tab. 7).

Tabella 7. *Composizione percentuale degli addetti e delle unità locali per attività economica in Puglia e province. Anno 2018.*

Attività economica	Puglia		Foggia		Barietta-Andria-Trani		Bari		Taranto		Brindisi		Lecce	
	Addetti	Unità locali	Addetti	Unità locali	Addetti	Unità locali	Addetti	Unità locali	Addetti	Unità locali	Addetti	Unità locali	Addetti	Unità locali
Industria	19,3	9,0	16,6	8,3	24,4	12,6	17,8	8,8	26,1	8,1	19,0	7,9	17,3	9,3
Costruzioni	8,9	10,5	8,6	10,4	7,6	9,4	9,0	10,2	7,1	9,0	9,6	11,5	10,6	12,0
Servizi, commercio, trasporto, alloggio e ristorazione	40,9	42,9	45,8	47,3	42,1	44,5	39,8	40,8	36,7	44,0	43,7	45,0	40,8	41,4
Servizi di informazione e comunicazione	2,1	1,7	1,0	1,2	1,1	1,3	3,1	2,2	1,7	1,6	1,1	1,4	1,9	1,6
Attività finanziarie e assicurative	2,5	2,3	2,7	2,2	2,2	2,1	2,6	2,3	2,1	2,4	2,0	2,1	2,6	2,2
Attività immobiliari	0,8	2,1	0,8	1,8	1,0	2,1	0,9	2,6	0,5	1,9	0,6	1,7	0,8	1,9
Attività professionali, scientifiche, tecniche ed amministrative	14,1	18,3	12,4	16,9	10,7	15,6	16,1	20,0	13,7	19,2	11,7	16,7	14,2	18,0
Istruzione, salute e servizi sociali	7,0	6,9	7,8	6,5	6,4	5,7	6,7	7,2	7,2	7,6	7,4	7,1	6,7	6,8
Altri servizi	4,5	6,2	4,3	5,4	4,5	6,8	4,0	6,0	5,0	6,3	4,8	6,7	5,1	6,7
Totale	100,0	100,0	100,0	100,0	100,0	100,0	100,0	100,0	100,0	100,0	100,0	100,0	100,0	100,0

Elaborazione su dati I.stat

4. Analisi delle dinamicità imprenditoriali sub regionali e comunali attraverso un'applicazione Shift-share

Alla luce dei differenti contesti territoriali delle strutture imprenditoriali, emersi con riferimento alle aree oggetto di studio, si è voluto approfondire il tema con una attenzione alle cause che potessero avere determinato tali risultanze. A tal fine si è proposta una particolare applicazione del modello di analisi Shift-share.

L'analisi Shift-share permette di "scomporre la variazione temporale di un fenomeno X (ad. es. numero imprese, UL, occupati o anche valore aggiunto delle imprese), le cui manifestazioni sono state rilevate contemporaneamente per settore di attività economica e per territorio, in modo tale da mettere in evidenza alcuni dei fattori [...] che possano averne influenzato lo sviluppo" (Biggeri *et al.* 2017).

In buona sostanza con l'analisi Shift-share la variazione assoluta dell'aggregato prescelto viene scomposta, con riferimento ad una sub-area territoriale, in tre componenti: tendenziale o della macro area (CM), strutturale (CS) e locale (CL).

In simboli si ha:

$$x_{ij,t} - x_{ij,0} = CM_{ij} + CS_{ij} + CL_{ij}$$

dove i è il macro-settore di attività, j è l'area geografica e t è il tempo.

La componente tendenziale misura la variazione che si sarebbe avuta nella sub-area territoriale in costanza di andamento complessivo del fenomeno nella macro area sovrastante

$$CM_{ij} = x_{ij,0} * r_{..}$$

dove $r_{..}$ è il tasso di variazione del fenomeno nell'intero paese o nella macro area di livello superiore.

La componente strutturale misura la parte di variazione attribuibile nella sub-area alla diversa composizione del tessuto produttivo (maggiore o minore presenza di attività "dinamiche")

$$CS_{ij} = x_{ij,0} *(r_i - r_{..})$$

dove r_i è il tasso di variazione del fenomeno nell'intero paese o nella macro area di livello superiore per il settore i .

Infine la componente residua, quella locale, esprime la quota di variazione espressione delle maggiori o minori specifiche capacità dinamiche degli operatori economici locali

$$CL_{ij} = x_{ij,0} *(r_{ij} - r_i)$$

dove r_{ij} è il tasso di variazione del fenomeno nel settore i e nell'area j .

Al fine di individuare le specificità delle imprese nelle province e nei comuni lucani e pugliesi, si è applicata l'analisi Shift-share al numero di addetti e alle UL delle imprese attive.

In particolare, date le variazioni assolute tra la fine e l'inizio del periodo in esame (2018 e 2012), per i due aggregati e per ciascuna provincia e comune, sono state individuate le tre componenti CM, CS e CL.

Nello specifico poi è importante evidenziare come, nel corso del presente lavoro, la componente tendenziale sia stata calcolata prendendo come riferimento la crescita media della regione e la componente strutturale considerando la suddivisione per macrosettore di attività economica riportata in appendice.

L'analisi sulla Basilicata (Tab. 8) rileva che se le due province avessero seguito il trend regionale, sia in termini di UL che di addetti, si sarebbe registrata una crescita maggiore a Potenza rispetto a Matera. Se il risultato era altamente prevedibile per quanto concerne le UL trattandosi di una analisi condotta sulle due sole province di cui una (Potenza) con un calo delle UL e l'altra (Matera) con un incremento, lo era meno con riferimento all'analisi condotta sugli addetti dove entrambe le province sono caratterizzate da un incremento del valore dell'aggregato osservato.

Tabella 8. *Variazione assoluta tra il 2018 e il 2012 e componenti dell'analisi Shift-share per la regione Basilicata.*

	Variazione assoluta 2018/2012	Componente tendenziale	Componente strutturale	Componente locale
Province	Unità locali			
Potenza	-66	121	-14	-173
Matera	248	61	14	173
	Addetti			
Potenza	6.870	7.724	-12	-841
Matera	4.549	3.695	12	841

Elaborazione su dati I.stat

A tal riguardo si può notare come la variazione assoluta di Matera risulti sostenuta da entrambe le componenti (strutturale e locale) sebbene sia decisamente più forte il contributo che ha dato alla crescita della provincia la componente locale, evidenziando quindi una particolare dinamicità del territorio unita ad una (seppur debole) maggior presenza di attività evidentemente trainanti.

Per la Puglia (Tab. 9), si evidenzia una differenziazione tra tre livelli di performance: le province più dinamiche sia in termini di UL che addetti (Bari e Barletta-Andria-Trani); la fascia intermedia leggermente sovraperformante per le UL e sottoperformante per gli addetti (Brindisi) e la fascia meno performante per entrambi gli aggregati (Taranto, Foggia e Lecce).

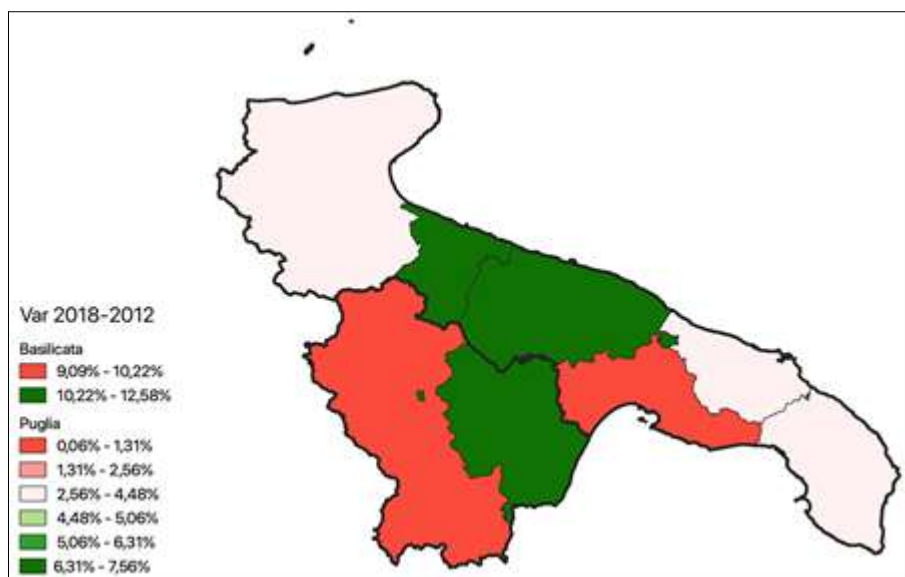
Tabella 9. *Variazione assoluta tra il 2018 e il 2012 e componenti dell'analisi Shift-share per la regione Puglia.*

Province	Variazione assoluta 2018/2012	Componente tendenziale	Componente strutturale	Componente locale
	Unità locali			
Foggia	-506	54	-86	-474
Bari	599	126	281	192
Taranto	-135	45	188	-368
Brindisi	107	34	-61	133
Lecce	-41	82	-192	69
Barletta-Andria-Trani	354	36	-130	448
Addetti				
Foggia	2.939	4.553	88	-1.703
Bari	19.631	13.006	640	5.984
Taranto	62	4.945	-113	-4.771
Brindisi	2.734	3.213	72	-552
Lecce	5.079	6.899	-269	-1.551
Barletta-Andria-Trani	5.337	3.164	-419	2.593

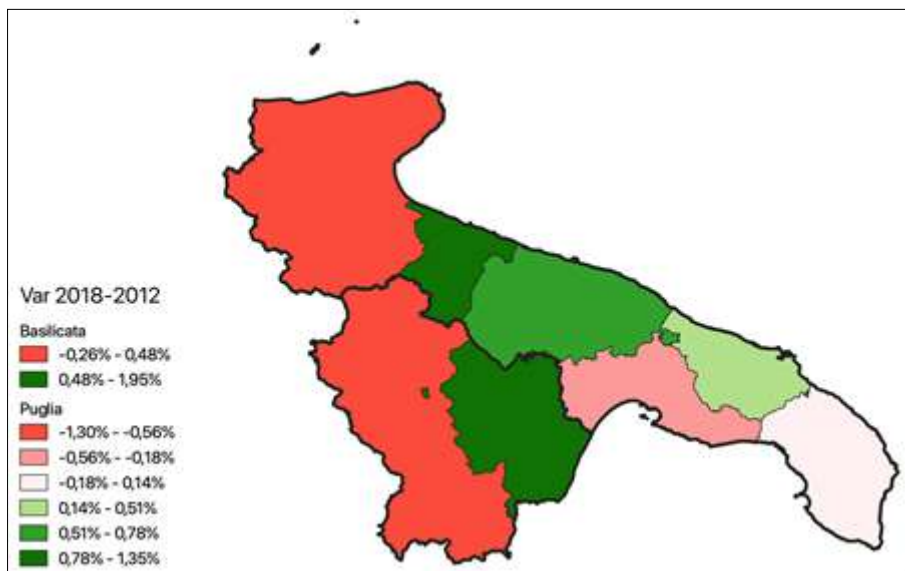
Elaborazione su dati I.stat

La crescita, di addetti e UL, delle province di Bari e Barletta-Andria-Trani è spiegata principalmente dalla componente locale particolarmente forte rispetto a quella evidenziatesi per le altre province, ovvero da una “indole imprenditoriale” dei due territori particolarmente marcata. Tale aspetto che accomuna le due province sembra essere ancora più marcato per la provincia di Barletta-Andria-Trani dove la struttura per attività economica delle imprese, e quindi la componente strutturale, incide negativamente, quasi a voler sottolineare una migliore capacità imprenditoriale caratterizzante l’area. Foggia e Taranto registrano una riduzione, nell’intero periodo, in termini di UL ed una crescita inferiore a quella media regionale in termini di addetti e tale risultanza appare visibilmente spiegata dal forte impatto negativo delle componenti locali.

Quanto sopra, è visibile in modo più semplice attraverso una lettura diversa, tramite i cartogrammi che seguono (Fig. 5 e Fig. 6) in cui sono riportate le province distinte in relazione alle variazioni relative percentuali di addetti e UL, evidenziando in verde le province che crescono di più della media regionale e in rosso quelle che crescono meno.

Figura 5. *Variazione percentuale degli addetti per provincia. Anno 2018 su anno 2012.*

Elaborazione su dati I.stat

Figura 6. *Variazione percentuale delle UL per provincia. Anno 2018 su anno 2012.*

Elaborazione su dati I.stat

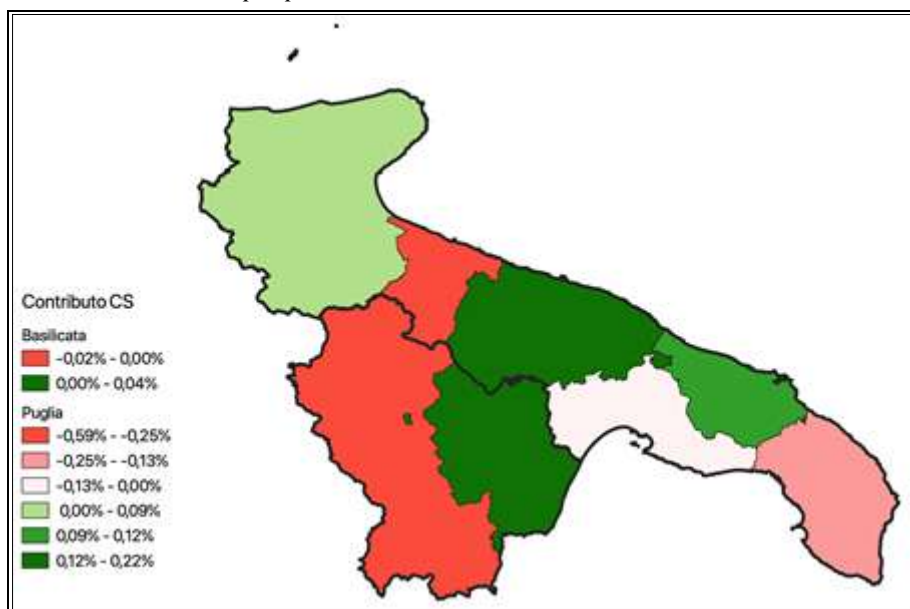
Nel periodo preso in esame, gli addetti della Basilicata crescono più del 10% a fronte di una crescita delle UL di appena lo 0,5%. È la provincia di Matera a crescere più rispetto alla media regionale sia in termini di addetti che di nuove UL.

In Puglia, la crescita nell'intero periodo è più contenuta rispetto alla Basilicata (3,5% gli addetti a fronte di un incremento quasi nullo delle UL) e sono le province di Bari e Barletta-Andria-Trani a registrare delle variazioni superiori alla media. Il minore incremento è registrato a Taranto e Foggia che hanno anche perso UL.

In modo analogo si è voluto rappresentare, per territori provinciali, il contributo percentuale delle componenti strutturali alla variazione degli addetti e delle UL, rappresentando in rosso le province in cui il contributo della componente è negativo ed in verde quelle in cui tale contributo è positivo (Fig. 7 e Fig. 8).

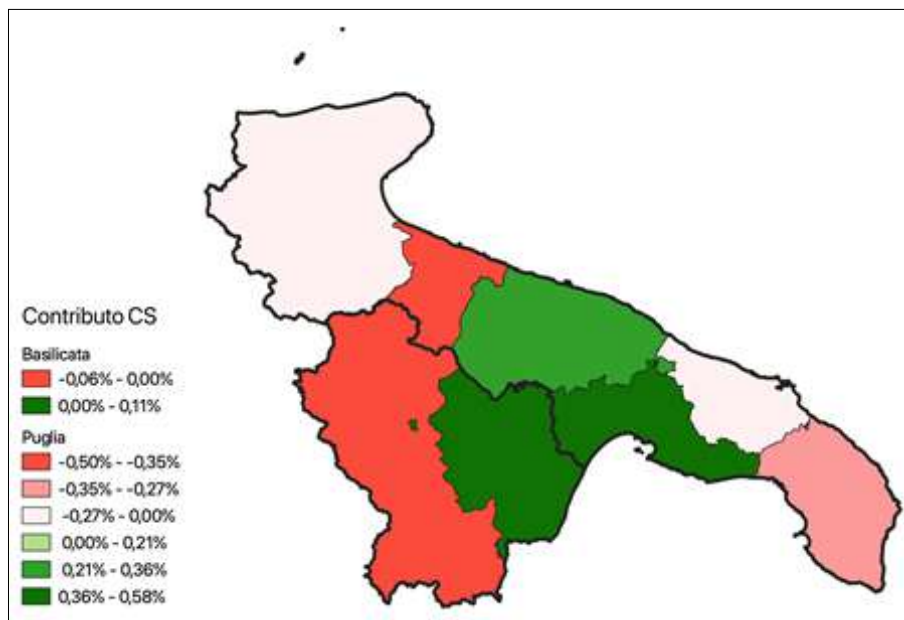
Si evince quindi che le particolari composizioni per attività economica hanno sostenuto la variazione positiva della provincia di Matera, per la Basilicata e, con riferimento alla Puglia, delle province di Bari e, in misura minore, Brindisi e Foggia per l'aggregato addetti mentre, risulta essere quasi nullo il contributo della componente strutturale, relativa alle UL, per Foggia.

Figura 7. Contributo percentuale alla variazione degli addetti della componente strutturale per provincia.



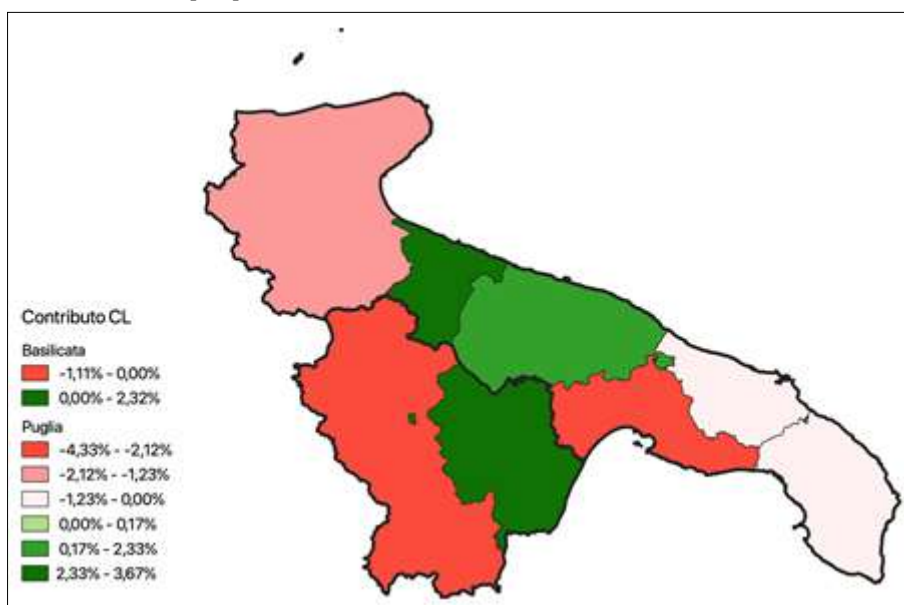
Elaborazione su dati I.stat

Figura 8. Contributo percentuale alla variazione delle UL della componente strutturale per provincia.



Elaborazione su dati I.stat

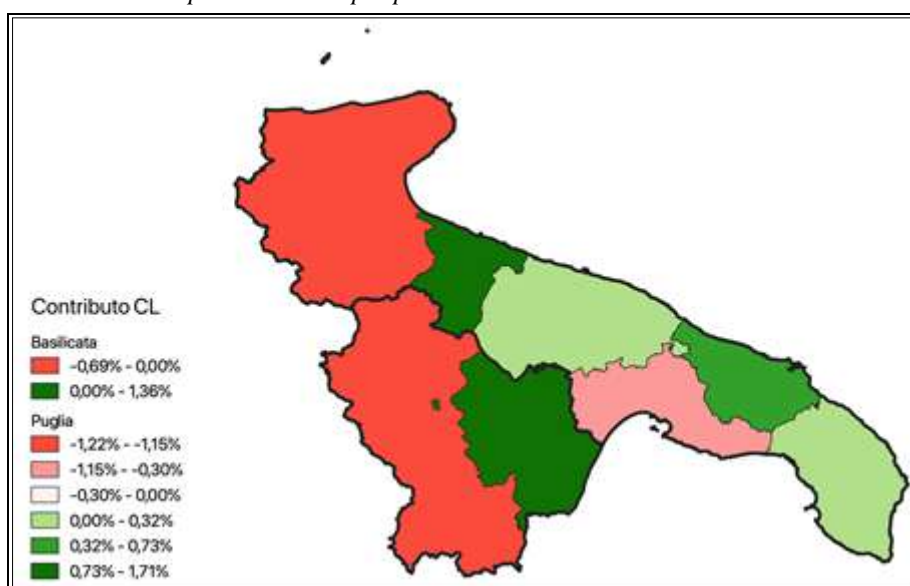
Figura 9. Contributo percentuale alla variazione degli addetti della componente locale per provincia.



Elaborazione su dati I.stat

La componente residuale o di dinamicità locale contribuisce in maniera positiva alla variazione degli addetti delle imprese localizzate nelle province di Matera, Bari, Barletta-Andria-Trani. Per le UL, tale componente continua a contribuire positivamente nelle tre province ed assume valore positivo anche nelle province di Brindisi e Lecce (Fig. 9 e Fig. 10).

Figura 10. *Contributo percentuale alla variazione delle unità locali della componente locale per provincia.*



Elaborazione su dati I.stat

Le risultanze così ottenute, per alcuni versi non di unica interpretazione, si sono volute condurre poi su base comunale al fine di individuare i comuni che hanno registrato un incremento superiore alla media regionale e quali siano le caratteristiche (di struttura o dinamicità locale) che hanno guidato tali incrementi.

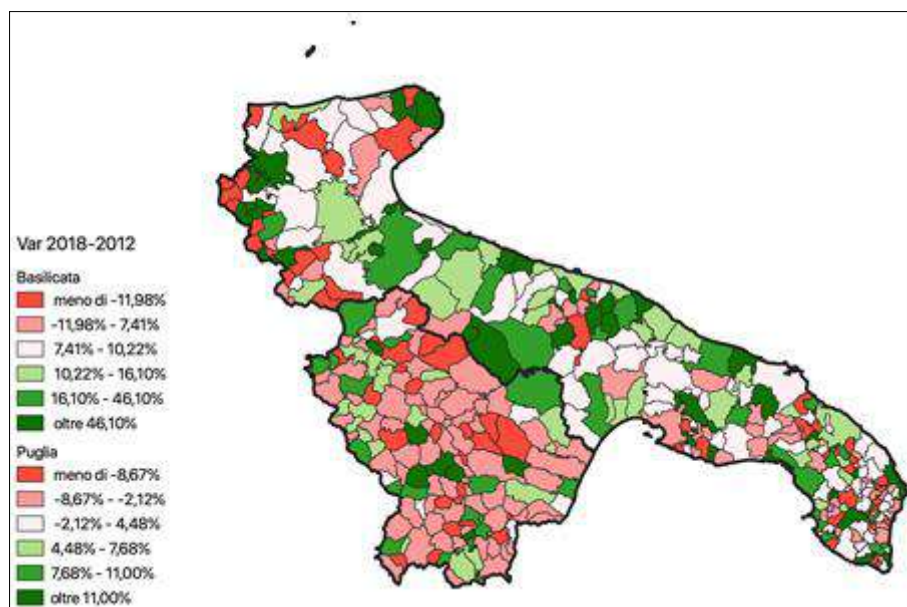
L'analisi su base comunale, se da un lato è potenzialmente molto utile nel definire le specificità territoriali, specialmente laddove, come nel presente caso, insistano profonde differenze locali, dall'altro rischia essere di non facile (se non addirittura in taluni casi fuorviante) lettura per comuni di piccole dimensioni in cui, per esempio, l'incremento in termini assoluti di poche UL, pesa molto in termini relativi, mostrando quindi variazioni a due cifre.

In generale, per la Basilicata e per la provincia di Matera, l'incremento relativo maggiore della media regionale, è registrato nella città di Matera, nei comuni

limitrofi e sulla litoranea, mentre per la provincia di Potenza è la Val d'Agri a registrare incrementi positivi per entrambi gli aggregati. Nel capoluogo lucano e nel lagonegrese, si registra invece una crescita più marcata delle UL che non si traduce in una crescita degli addetti.

In Puglia e in particolare per la provincia di Bari, il capoluogo di regione, il sud est barese e l'alta murgia registrano incrementi per entrambi gli aggregati. In generale, per le UL, molte aree presentano incrementi maggiori rispetto alla media, sebbene a tali incrementi non corrisponda la stessa variazione positiva in termini di addetti (Fig. 11 e Fig. 12).

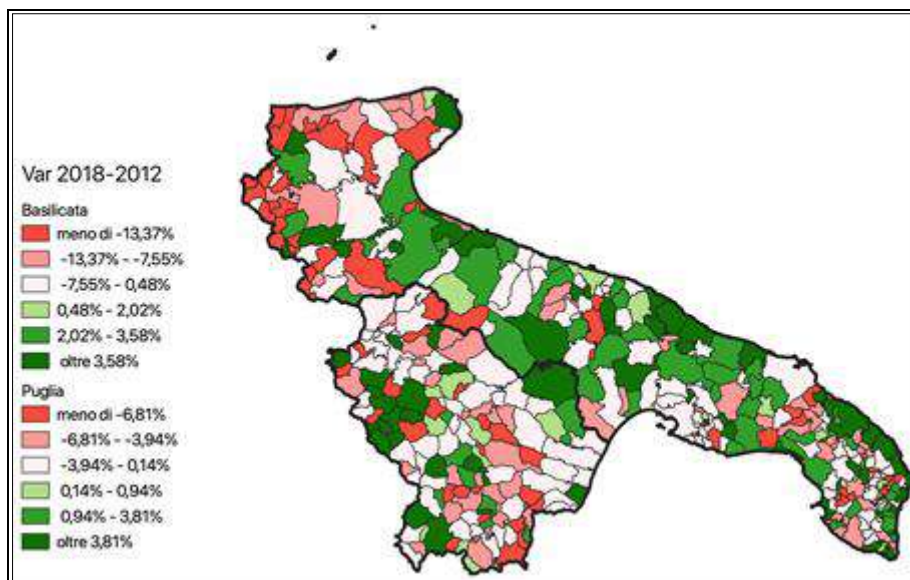
Figura 11. *Variazione percentuale degli addetti per comune. Anno 2018 su anno 2012.*



Elaborazione su dati I.stat

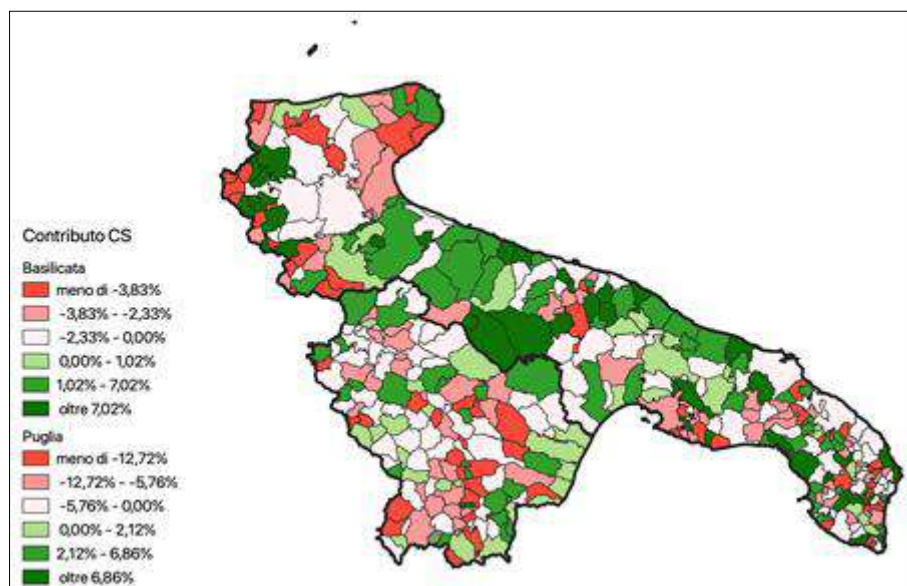
In generale si evince come, mentre con riferimento agli addetti appare esserci una sorta di epicentro virtuoso in Puglia prevalentemente diffuso lungo le coste e nella zona murgiana, tale andamento positivo tende a scemare allontanandosi da questa area e giungendo in Basilicata. Con riferimento alle UL si evidenziano zone di vivacità meno concentrate in Basilicata.

Figura 12. *Variazione percentuale delle unità locali per comune. Anno 2018 su anno 2012.*



Elaborazione su dati I.stat

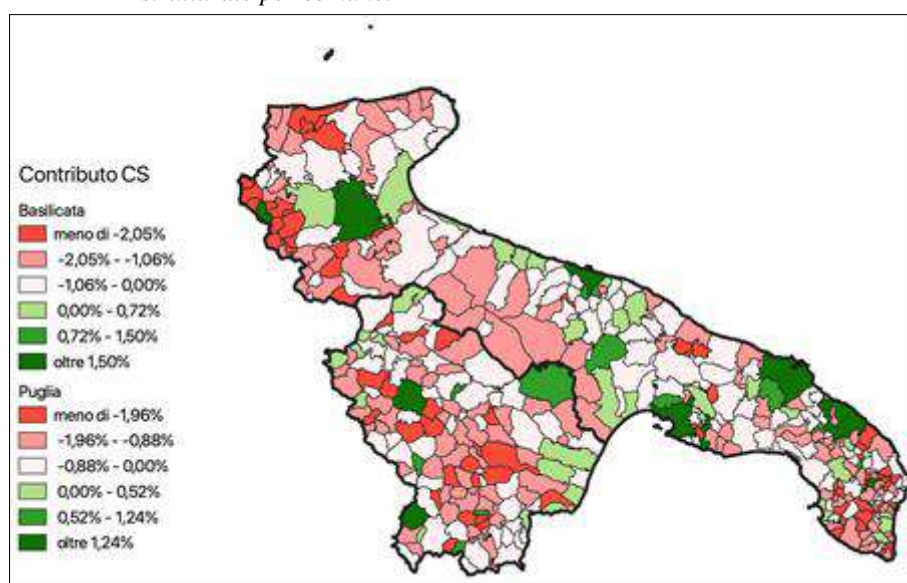
Figura 13. *Contributo percentuale alla variazione degli addetti della componente strutturale per comune.*



Elaborazione su dati I.stat

In generale, una importante spinta all'incremento dei due aggregati, nell'intero periodo in esame, è data dalla presenza sul territorio di attività economiche che hanno registrato un incremento maggiore rispetto alla media. La componente strutturale, presentata, a livello comunale, in termini di contributo percentuale alla variazione (Fig. 13 e Fig. 14), sintetizza tale aspetto. In Basilicata in corrispondenza di Matera, di Potenza, della costa ionica e del melfese, aree a vocazione maggiormente industriale, la componente strutturale ha contribuito in modo positivo alla variazione dei due aggregati.

Figura 14. *Contributo percentuale alla variazione delle UL della componente strutturale per comune.*



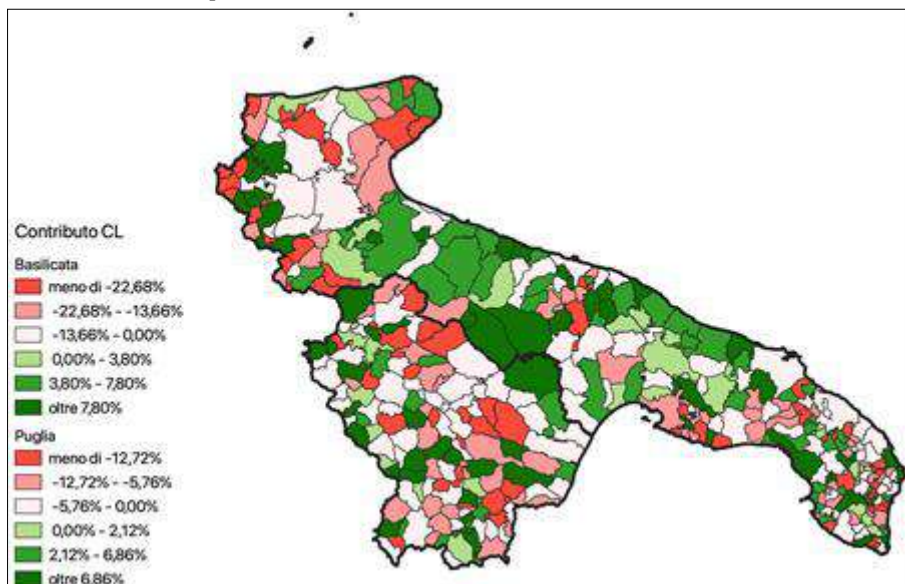
Elaborazione su dati I.stat

In Puglia, in termini di addetti, vi è un buon impatto della componente strutturale nelle stesse aree che hanno registrato un maggiore incremento dell'aggregato ma, la stessa cosa non può dirsi in termini di UL. Molte aree della Puglia, pur presentando una crescita, rispetto alla media regionale, registrano una componente strutturale negativa o debolmente positiva. La crescita in termini di UL non sarebbe quindi prevalentemente dettata dalla composizione per attività economica dei due aggregati nei diversi comuni.

Esaminando, infine, il contributo percentuale alla variazione dei due aggregati della componente locale per comune (Fig. 15 e Fig. 16), in Basilicata si nota come tale componente, pur pesando di più rispetto a quella strutturale, ricalchi lo stesso andamento per gli addetti. Nelle aree in cui vi sono attività economiche che hanno

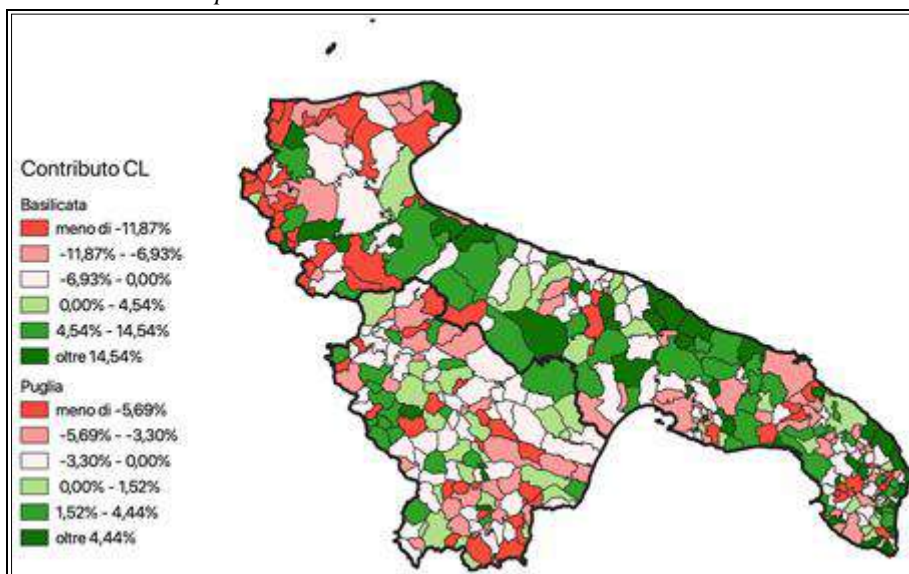
supportato la crescita degli addetti, vi sono stati anche fattori locali di dinamicità che hanno contribuito all'incremento positivo.

Figura 15. *Contributo percentuale alla variazione degli addetti della componente locale per comune.*



Elaborazione su dati I.stat

Figura 16. *Contributo percentuale alla variazione delle UL della componente locale per comune.*



Elaborazione su dati I.stat

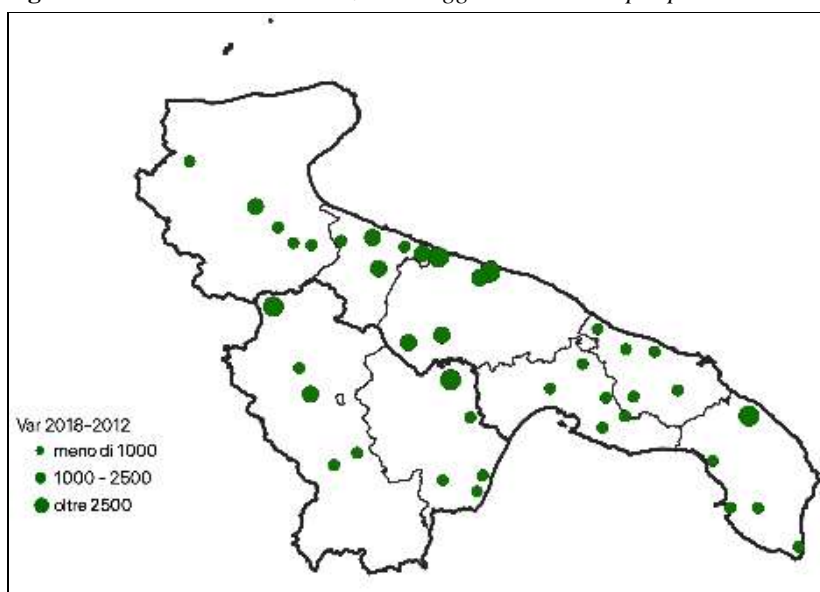
In termini di UL si registra un andamento della componente locale meno simile rispetto alla componente strutturale. Vi sono aree, principalmente intorno al capoluogo di regione, in cui comunque si è registrata una spinta imprenditoriale locale, sebbene per attività a bassa redditività.

La stessa cosa dicasi per la Puglia dove, in presenza di una componente strutturale negativa, si registra una componente locale positiva per le UL che non si traduce comunque in un incremento della componente locale calcolata sugli addetti.

Al fine di entrare in modo più puntuale nelle specificità comunali, si riportano (Fig. 17 e Fig. 18), per ciascuna provincia e ciascun aggregato, i cinque comuni che hanno registrato la variazione assoluta maggiore degli addetti e delle UL.

Da tale analisi si evince chiaramente come, in riferimento ad entrambi gli aggregati, esista un fenomeno “gravitazionale” in virtù del quale i comuni con migliori performance tendono ad essere vicini tra loro.

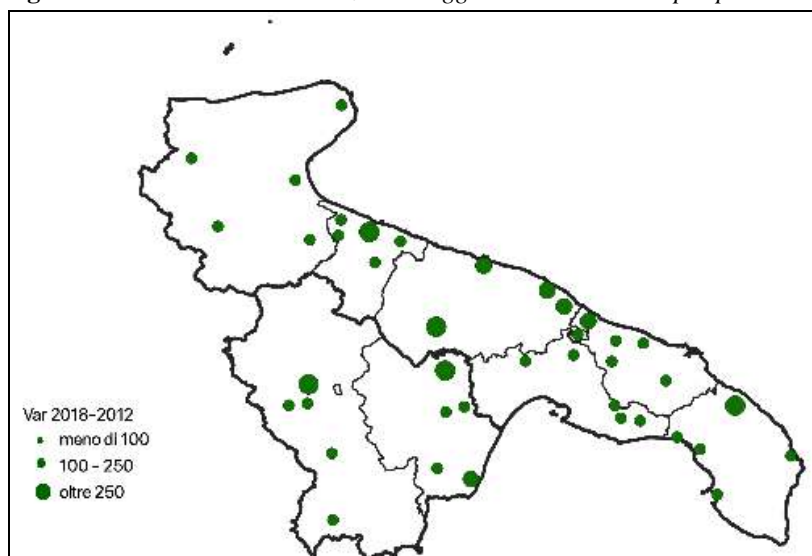
Figura 17. Comuni con la variazione maggiore di addetti per provincia.



Elaborazione su dati I.stat

In Basilicata, l’incremento maggiore in termini di UL è registrato nei due capoluoghi di provincia, ma Potenza ha una struttura economica che comunque è cresciuta maggiormente nel periodo in esame. Matera invece deve la sua crescita a fattori “locali”: ricordiamo che si preparava all’evento Matera 2019 che ha dato grosso slancio dal punto di vista economico e occupazionale alla città.

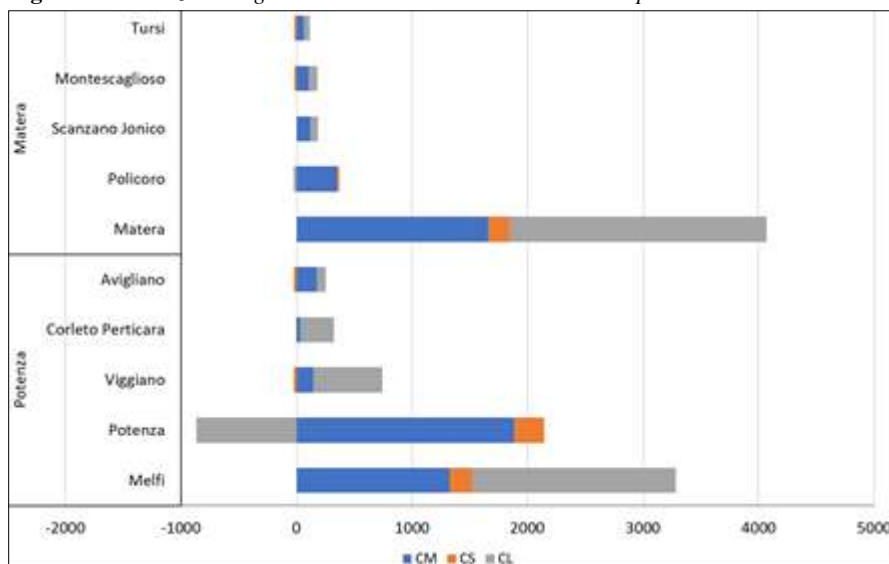
Figura 18. Comuni con la variazione maggiore di unità locali per provincia.



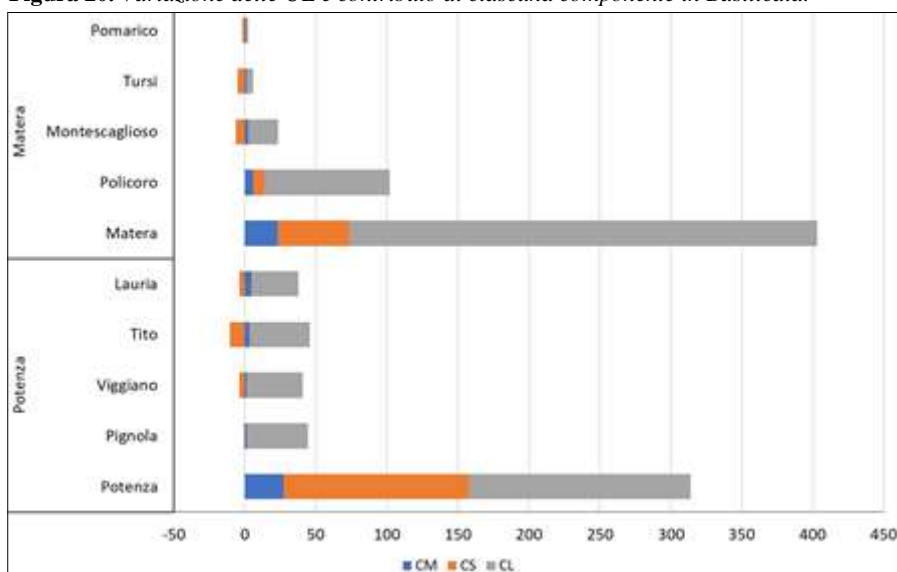
Elaborazione su dati I.stat

In termini di occupati, il primo posto è assegnato, per la provincia di Potenza, a Melfi, che registra una componente locale positiva, a differenza di Potenza, che pur occupando il secondo posto presenta una componente locale negativa. Anche in termini di addetti, Matera si conferma al primo posto, seguita da Policoro (Fig. 19 e Fig. 20).

Figura 19. Variazione degli addetti e contributo di ciascuna componente in Basilicata.

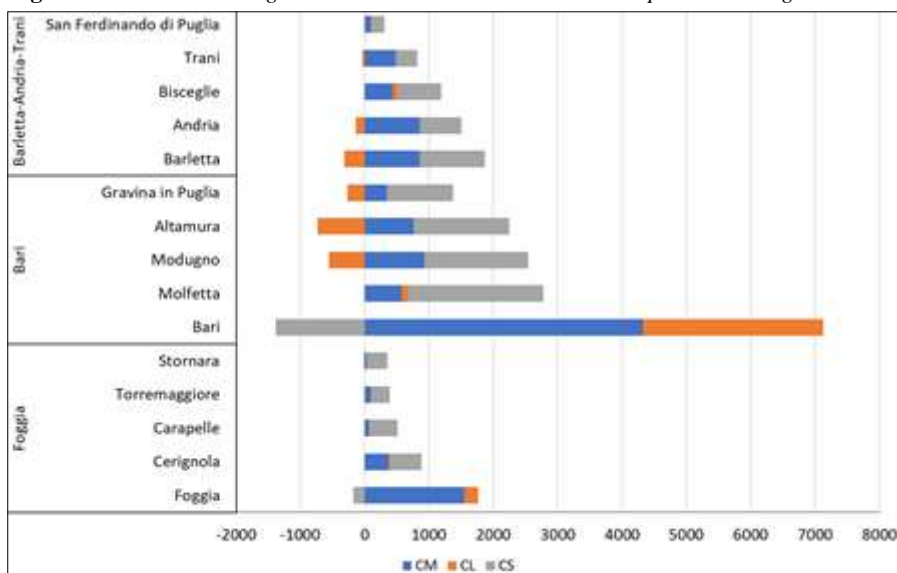


Elaborazione su dati I.stat

Figura 20. *Variazione delle UL e contributo di ciascuna componente in Basilicata.*

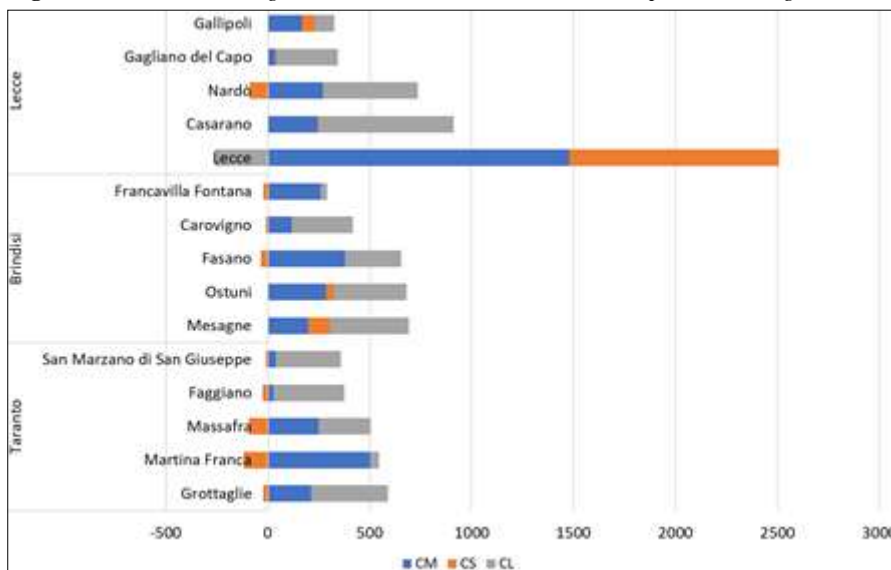
Elaborazione su dati I.stat

Per la Puglia, in termini di addetti, tra i capoluoghi di provincia, solo Brindisi e Taranto non si posizionano nei primi posti per incremento; nella provincia di Brindisi, ai primi posti, si posizionano Mesagne, Ostuni e Fasano; nella provincia di Taranto registrano la crescita più alta Grottaglie, Martina Franca e Massafra.

Figura 21a. *Variazione degli addetti e contributo di ciascuna componente in Puglia.*

Elaborazione su dati I.stat

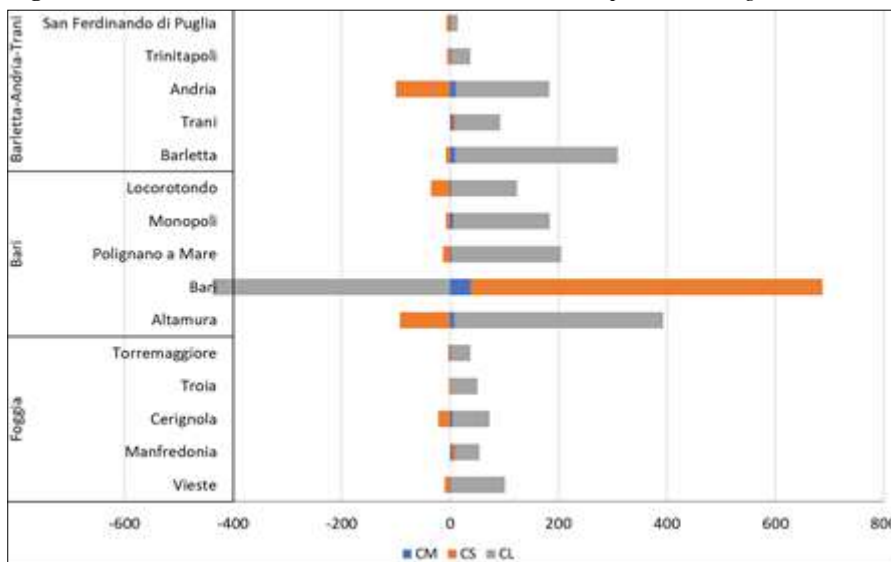
Figura 21b. *Variazione degli addetti e contributo di ciascuna componente in Puglia.*



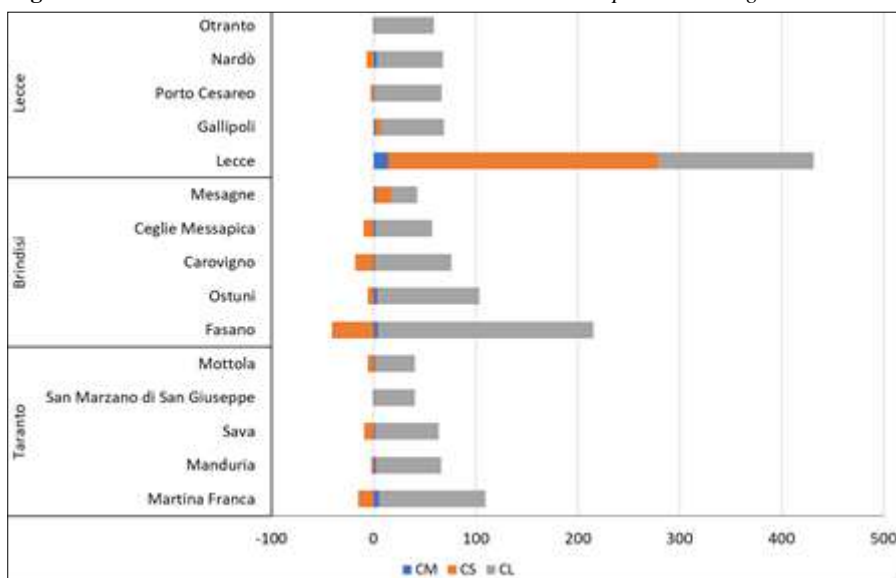
Elaborazione su dati I.stat

Nella provincia di Barletta-Andria-Trani le tre città sono tra quelle con crescita maggiore sia per addetti che UL; Bari, nonostante abbia una componente strutturale pesantemente negativa, si posiziona al primo posto per gli addetti ed al secondo per le UL, preceduta da Altamura (Figg. 21 e Figg. 22).

Figura 22a. *Variazione delle UL e contributo di ciascuna componente in Puglia.*



Elaborazione su dati I.stat

Figura 22b. *Variazione delle UL e contributo di ciascuna componente in Puglia.*

Elaborazione su dati I.stat

5. Alcune considerazioni conclusive

Lo studio condotto e qui illustrato evidenzia come i dati derivanti dal Registro ASIA-UL, siano uno strumento decisamente idoneo per applicazioni con analisi metodologiche non particolarmente complesse, quali le varie declinazioni applicative Shift-share.

Tali applicazioni permettono l'individuazione e la definizione di specificità territoriali sub-provinciali caratterizzate da particolari significatività in termini di dinamicità altrimenti difficilmente rilevabili, se non con l'utilizzo di fonti informative non Sistan, quindi non aventi il carattere di "ufficialità" e con differente significatività.

L'aspetto più interessante lo si evince non solamente nella idoneità emersa da parte del Registro stesso, quanto nella idoneità del particolare dettaglio informativo dei dati utilizzati.

Infatti, i dati elaborati nel presente lavoro, si ricorda, non sono i file dei microdati, per antonomasia atti ad avere una maggiore malleabilità di elaborazione e quindi un maggior grado di precisione nelle analisi, bensì i dati aggregati così come presenti nella banca dati I.stat.

Ed è proprio l'unione tra estrema agilità e universalità di offerta informativa Sistan (accessibile a tutti senza alcun bisogno di autorizzazioni e motivazioni per il

loro utilizzo), e bontà del modello di analisi Shift-share emersa nel lavoro, a destare interesse.

L'unico limite riscontrato consta, forse, nella assenza dei dati relativi al valore aggiunto (non presente nei Registri ASIA e ASIA-UL) che avrebbero apportato una maggiore qualità di analisi e di interpretazione, consentendo di introdurre anche valutazioni basate sulla produttività delle imprese delle diverse aree territoriali.

Riferimenti bibliografici

- Bianchino, A.; Carucci, A.M.M.; Vannella, G. (2018). *Il nuovo sistema statistico "a misura di comune": una applicazione sul territorio lucano*. Metodi e analisi statistiche, Università degli Studi di Bari Aldo Moro
- Biggeri, L.; Bini, M.; Coli, A.; Grassini, L.; Maltagliati, M. (2017). *Statistica per le decisioni aziendali*. Pearson Italia
- Boselli, C.; Brunetti, S.; Cammarota, M.; De Giorgi, V.; D'Urzo, A.; Ricci, M.; Pazzini, R.; Seri, G.; Siesto, G.; Virgili, L. (2016). *Il processo di diffusione dei dati delle statistiche strutturali sulle imprese (Frame-SBS): aspetti normativi e metodologici connessi all'ampliamento del dettaglio informativo*. Istat working papers n. 14. 2016
- Carucci, A.M.M.; Vannella, G. (2016). *Sull'integrazione tra fonti amministrative e statistiche per le imprese*. Metodi e analisi statistiche, Università degli Studi di Bari Aldo Moro
- Carucci, A.M.M.; Vannella, G. (2019). *Il Frame-SBS territoriale: struttura e potenzialità interpretative tramite analisi Shift-share*. Metodi e analisi statistiche, Università degli Studi di Bari Aldo Moro
- Faramondi, A.; De Giorgi, V.; De Francesco, D.; Di Manno, R.; Lombardi, S.; Nardecchia, R.; Sanzo, R.; Tomeo, V.; Trinca E. (2018). *La stima del valore aggiunto a livello territoriale: il nuovo registro statistico "Frame-SBS Territoriale"*, Atti della XXXIX Conferenza italiana di scienze regionali
- Seri, G.; Ichim, D.; Luchetti, F.; Costa, S.; Nurra, A.; Mastrostefano, V.; Salamone, S.; Pascucci, C.; Orsini, D. (2016). *Integrazione del Frame con altre indagini e fonti amministrative ai fini della produzione di indicatori complessi*. Istat working papers n. 17. 2016

APPENDICE

Raccordo macrosettori di attività economiche e sezioni Ateco 2007.

Macrosettore di attività economica	Ateco 2007
Agricoltura	A - AGRICOLTURA, SILVICOLTURA E PESCA
Industria	B - ESTRAZIONE DI MINERALI DA CAVE E MINIERE C - ATTIVITÀ MANIFATTURIERE D - FORNITURA DI ENERGIA ELETTRICA, GAS, VAPORE E ARIA CONDIZIONATA E - FORNITURA DI ACQUA; RETI FOGNARIE, ATTIVITÀ DI GESTIONE DEI RIFIUTI E RISANAMENTO
Costruzioni	F - COSTRUZIONI
Servizi, commercio, trasporto, alloggio e ristorazione	G - COMMERCIO ALL'INGROSSO E AL DETTAGLIO; RIPARAZIONE DI AUTOVEICOLI E MOTOCICLI H - TRASPORTO E MAGAZZINAGGIO I - ATTIVITÀ DEI SERVIZI DI ALLOGGIO E DI RISTORAZIONE
Servizi di informazione e comunicazione	J - SERVIZI DI INFORMAZIONE E COMUNICAZIONE
Attività finanziarie e assicurative	K - ATTIVITÀ FINANZIARIE E ASSICURATIVE
Attività immobiliari	L - ATTIVITÀ IMMOBILIARI
Attività professionali, scientifiche, tecniche ed amministrative	M - ATTIVITÀ PROFESSIONALI, SCIENTIFICHE E TECNICHE N - NOLEGGIO, AGENZIE DI VIAGGIO, SERVIZI DI SUPPORTO ALLE IMPRESE
Istruzione, salute e servizi sociali	P - ISTRUZIONE Q - SANITÀ E ASSISTENZA SOCIALE
Altri servizi	R - ATTIVITÀ ARTISTICHE, SPORTIVE, DI INTRATTENIMENTO E DIVERTIMENTO S - ALTRE ATTIVITÀ DI SERVIZI

Data Envelopment Analysis: application of the efficient frontier on the financial field in the European and the American scenarios

Crescenzo Gallo^{1*}, Alessandro Rinaldi²

¹ Department of Clinical and Experimental Medicine, University of Foggia

² Department of Economics and Finance University of Bari A. Moro

Abstract: The topic of efficiency analysis of independent organisational units has long been a major challenge for researchers and scientists around the world. In fact, in the current economic context of increasing competitiveness and dynamism, it is essential to know both the degree of efficiency compared to competitors, and the efficiency of the various internal operating units or individual dependent units. This technique, called Data Envelopment Analysis (DEA)¹, can be applied in different contexts ranging from the economic to the financial one, essentially based on the identification of the best targets with which to compare the performance of various units. These can be banks, companies, insurance companies or securities portfolios, so that at the same time they can represent, evaluate and improve the positions and performance of each of them. In the light of the above, the present work uses the Data Envelopment Analysis methodology used in the financial sector, in order to identify the optimal and efficient frontier of securities relating to the European and American areas which, on the basis of Harry Markowitz's theory, represent the whole of those portfolios (known as *dominant*) which, for the same yield, are the least risky or which, for the same risk, are the most profitable.

In order to achieve this objective, it is essential that the securities in the portfolio under study are not perfectly correlated. By comparing the frontiers obtained, through the analysis of the volatility and expected profit of each portfolio, it is possible to arrive at a knowledge of the market that is more convenient and advantageous in terms of risk/return for any investor.

Keywords: DEA; Efficient frontier; VAR.

* Corresponding author: crescenzo.gallo@unifg.it.

This work is the result of a common project: however, A. Rinaldi provided for the preparation of paragraph 2 and 4, C. Gallo drafted paragraph 1 and 3.

¹PREMACHANDRA I.M., CHEN Y., WATSON J., (2011), "DEA as a tool for predicting corporate failure and success: A case of bankruptcy assessment".

1. Introduction

In the last decade, the study of the efficiency of decision-making units has been a highly debated topic in order to improve their performance in a purely economic and financial context.

DEA (Data Envelopment Analysis)² was first introduced by Charnes, Cooper and Rhodes³ in 1978 with the publication of “Measuring the efficiency of decision-making units”. In this work the authors aimed to calculate the relative efficiency⁴ of similar production units (Decision Making Units, DMU) characterized by a system of input and output, that is, separate units within a management organization chart. In this perspective, starting from Farrel’s efficiency index⁵, the authors managed to transform the analysis he had conducted on a single output into a linear programming problem so as to apply this theory to a multitude of outputs. To this end, the objective of their work was to extend Farrel’s idea of the possibility of being able to estimate technical efficiency through a production function capable of outlining the maximum output level achievable for a given input level in a technically efficient manner. From these writings and publications, the economist Markowitz elaborated a theory based on a series of techniques, methods and procedures that generate market supply and demand according to the risk/return ratio, with the intent of being able to build an efficient portfolio of securities in order to minimize the risk and maximize the overall return that would result from them.

Hence the concept of *efficient frontier*, on which the idea of Markowitz is based, composed precisely of all those securities portfolios placed along a curve that for a given yield have less risk and for a given degree of risk have more yield. For this to be the case, however, the securities making up the portfolio must not be correlated with each other, i.e. there must be no correlation of any kind between them.

In this article, the final aim is to illustrate and compare, within the DEA, the efficient borders of securities portfolios belonging to the European and American scenarios.

²CHARNES A., COOPER W.W., LEWIN A., SEINFORD L. (1994), “Data Envelopment Analysis Theory, Methodology and Applications”.

³CHARNES A., COOPER W., RHODES E. (1978), “Measuring the Efficiency of Decision Making Units, European Journal of Operational Research”.

⁴BISCEGLIA M., (2015), “Il controllo di gestione”. Studies in memory of "Carlo Cecchi".

⁵FARRELL M. (1957), “The Measurement of Productive Efficiency”.

2. The efficient frontier from the Markowitz's perspective

The theory on which the concept of efficient frontier according to the Markowitz's perspective is based, concerns a combination of asset classes that maximizes the expected return given a certain level of risk or, equivalently, minimizes the risk given a level of expected return. The union of these optimal portfolios graphically traces a curve called the *frontier*, which lies in a two-dimensional space (risk-expected return). The portfolios that make up the frontier are efficient, i.e. there is no prior combination that can overcome them in terms of trade-off risk-return.

All portfolios lying below the border are sub-optimal, as it is possible to find one that offers a higher expected return for the same risk or one that has a lower risk for the same expected return. The fundamental cornerstones of Markowitz's portfolio theory can therefore be summarised in 5 key points:

- Investors want to maximize the final wealth by trying to minimize the level of risk;
- The period for which the investment is intended to be made must be predefined by a single interval defined by the range: $(P_t; P_{t+T})$;
- Transaction costs and taxes must be zero;
- The expected return and the average squared deviation or standard deviation (which represents volatility and in this case the risk), are the only parameters on which the final decision must converge;
- The market in which the investment is made must be in perfect competition.

The expected return (or mathematical expectation) of a financial asset is given by the formula:

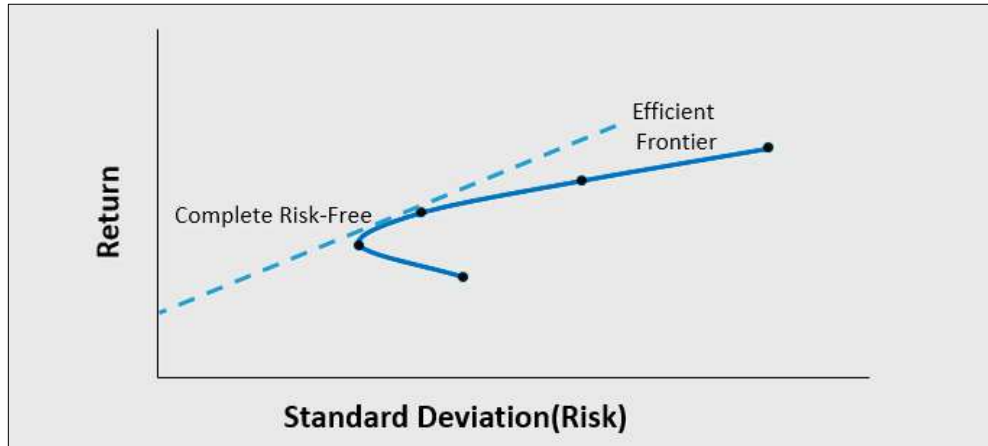
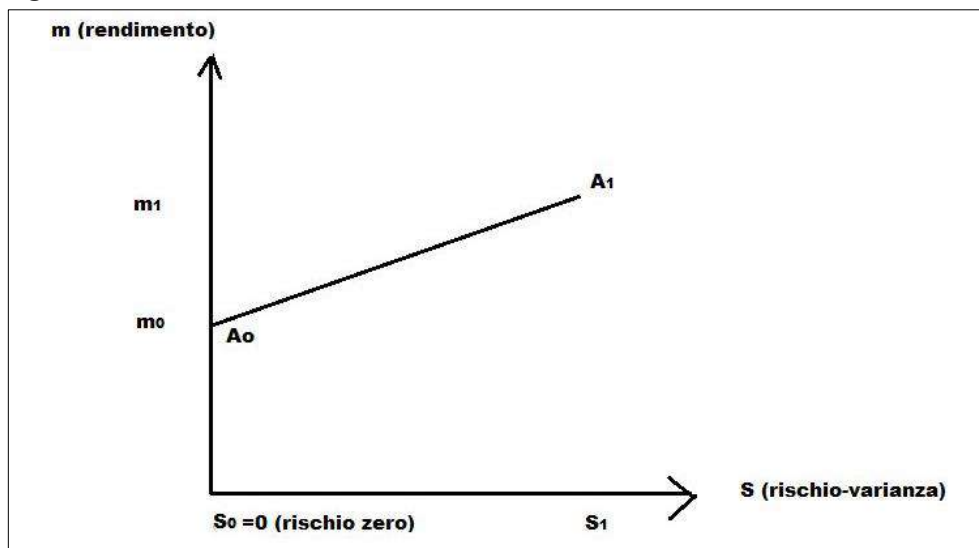
$$E(X) = \mu_x = \sum_{i=1}^n x_i p(x_i) \quad (1.1)$$

where x_i represents the set of values that a random variable can express, while $p(x_i)$ represents the probability that the variable can assume. In the case in which we are in front of a random variable constituted by infinite measurements we would have:

$$E(X) = \mu_x = \int_{-\infty}^{+\infty} x f_x(x) dx \quad (1.2)$$

instead of the expected risk, which can be defined as the "degree of uncertainty that the market expresses on the actual achievement of expected returns" (Fig.1).

Before talking about an efficient frontier, a premise must be made by first considering the combination of two securities $A_0 = (\sigma_0, m_0)$ and $A_1 = (\sigma_1, m_1)$, where σ_0, σ_1 indicate the risk/variance and m_0, m_1 represent the yield of the two securities (Fig.2).

Figure 1 – An Example of Efficient Frontier from the market point of view**Figure 2-** Efficient frontier between two securities A0 and A1

The security A_0 is risk free or *risk-free*, while A_1 possesses a high risk. The segment connecting A_0 with A_1 represents the efficient frontier, with:

$$1) \text{ return } m = (x_1 m_1) + (x_2 m_2)$$

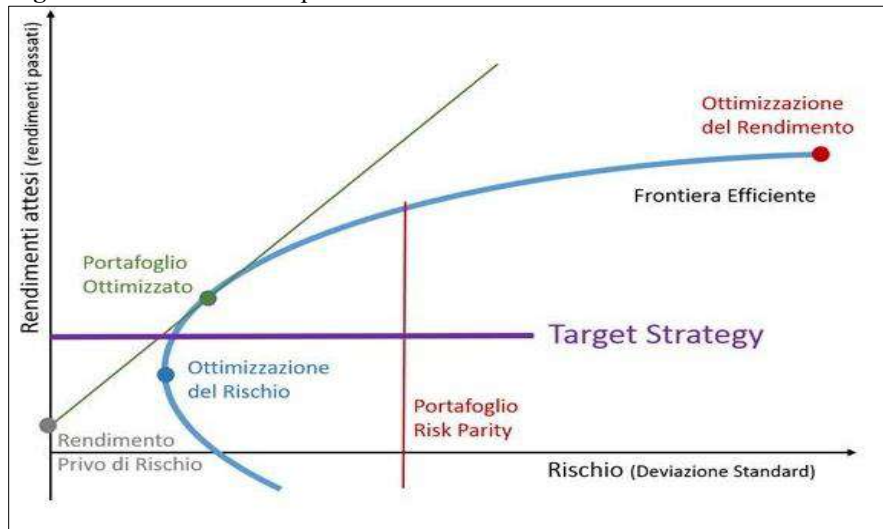
$$2) \text{ variance } \sigma^2 = x_1^2 \sigma_1^2 + x_2^2 \sigma_2^2 + 2x_1 x_2 \sigma_{1,2} \rho_{1,2}$$

where $\sigma_{1,2}$ and $\rho_{1,2}$ indicate respectively the covariance and the correlation between the two securities.

Consider n securities A_j (Fig.3) where:

- m_j representing the expected return of the j -th security;
- σ_j representing the variance of the j -th security .

Figure 3- Mean / Variance plane and Efficient Frontier



By combining the n amounts of capital to be distributed over the securities, the total return is:

$$m = \sum_{j=1}^n x_j m_j \tag{1.3}$$

while the overall variance will be given:

$$\sigma^2 = \sum_{j=1}^n x_j^2 \sigma_j^2 + 2 \sum_{h < k} x_h x_k \sigma_{h,k} \tag{1.4}$$

This formula is needed to construct the matrix of variances-covariances

$$\sum_x = \begin{pmatrix} \sigma_{x_1}^2 & \sigma_{x_1,x_2} & \dots & \sigma_{x_1,x_i} & \dots & \sigma_{x_1,x_n} \\ \sigma_{x_2,x_1} & \sigma_{x_2}^2 & \dots & \sigma_{x_2,x_i} & \dots & \sigma_{x_2,x_n} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \sigma_{x_i,x_1} & \sigma_{x_i,x_2} & \dots & \sigma_{x_i}^2 & \dots & \sigma_{x_i,x_n}^2 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \sigma_{x_n,x_1} & \sigma_{x_n,x_2} & \dots & \sigma_{x_n,x_i}^2 & \dots & \sigma_{x_n}^2 \end{pmatrix}$$

The variance can also be written as:

$$\sigma^2 = \sum_{h,k=1}^n x_h x_k \sigma_{h,k} \quad (1.5)$$

Assuming $\sum_{j=1}^n x_j = 1 \rightarrow$ (*first linear constraint*) it is possible to sell short⁶ the

Stocks (*short selling*) and, employing the method of Lagrange multipliers, it will be possible to find and to define the *efficient frontier*.

A Lagrangian function will be of the type $f(x, y)$ with the constraint $g(x, y) = c$. Let's consider a real parameter λ to construct the Lagrangian function

$$L(x, y, \lambda) = f(x, y) - \lambda(g(x, y) - c) \quad (1.6)$$

whose gradient will be

$$\nabla L(x, y, \lambda) = 0 \Rightarrow \begin{cases} L'_x(x, y, \lambda) = f'_x(x, y) - \lambda g'_x(x, y) \\ L'_y(x, y, \lambda) = f'_y(x, y) - \lambda g'_y(x, y) \\ L'_\lambda(x, y, \lambda) = -(g(x, y) - c) \end{cases} \quad (1.7)$$

If we instead assume $\sum_{j=1}^n x_j \geq 0 \rightarrow$ (*second linear constraint*) it is not possible to sell short.

Let's take a frontier portfolio

$$F(m) = x_i(m), \quad x_i = x_1, x_2, \dots, x_n \quad (1.8)$$

with performance m and components $x_i(m)$, where $x_i(m)$ is a linear equation dependent on m of type:

$$x_i(m) = a_i m + b_i \quad (1.9)$$

and

$$\sum_{i=1}^n x_i(m) = \sum_{i=1}^n a_i m + b_i = 1 \quad (1.10)$$

From the above it follows that

$$\sum_{i=1}^n a_i m + \sum_{i=1}^n b_i = m \sum_{i=1}^n a_i + \sum_{i=1}^n b_i = 1 \quad (1.11)$$

⁶The sale of shares in the open (in English short selling or short) is a financial transaction involving the sale of equity securities without having the properties, hoping to repurchase them later at the time to deliver them to the buyer at a lower price, thus realizing a profit; a variant of this type of operation is to borrow securities to be someone and to return later in the hope that prices may have fallen.

By placing the constraints

$$\sum_{i=1}^n a_i = 0 \quad e \quad \sum_{i=1}^n b_i = 1 \quad (1.12)$$

if we take two frontier portfolios P_1 and P_2 , the composite portfolio will always be a frontier portfolio. In fact, considering the two portfolios:

- $P_1 = x_i$ with performance $F(m_1) = x_i = m_1$
- $P_2 = y_i$ with performance $F(m_2) = y_i = m_2$

the composite portfolio will be of the type:

$$\begin{aligned} F &= kP_1 + (1-k)P_2 \Rightarrow F = kF(m_1) + (1-k)F(m_2) \Rightarrow \\ &\Rightarrow F = km_1 + (1-k)m_2 \end{aligned} \quad (13)$$

By placing the constraints:

$$m_1 = a_i m_1 + b_i \quad e \quad m_2 = a_i m_2 + b_i \quad (1.14)$$

the composition of the portfolio will be:

$$\begin{aligned} \sum_{i=1}^n k(a_i m_1 + b_i) + \sum_{i=1}^n (1-k)(a_i m_2 + b_i) &\Rightarrow k(a_i m_1 + b_i) + (1-k)(a_i m_2 + b_i) \Rightarrow \\ &\Rightarrow ka_i m_1 + kb_i + (1-k)a_i m_2 + (1-k)b_i \end{aligned}$$

Finally setting out a_i you get

$$a_i (km_1 + (1-k)m_2) = m_3 \quad (1.15)$$

which will represent the overall return on the frontier portfolio given by the combination of two frontier portfolios.

If we also consider a frontier portfolio such as $F_0 = (y_j)$ the covariance between the portfolio F_0 and the portfolio $F(m) = x_i(m)$ will be given by:

$$Cov(F_0, F(m)) = \sum_{i=1}^n y_j (a_i m + b_i) \sigma_{i,j} \quad (1.16)$$

from which you will get:

$$Cov(F_0, F(m)) = \sum_{i=1}^n y_j a_i m \sigma_{i,j} + \sum_{i=1}^n y_j b_i \sigma_{i,j} = m \sum_{i=1}^n y_j a_i \sigma_{i,j} + \sum_{i=1}^n y_j b_i \sigma_{i,j} \quad (1.17)$$

By putting:

$$\alpha = \sum_{i=1}^n y_j a_i \sigma_{i,j} \quad \text{is} \quad \beta = \sum_{i=1}^n y_j b_i \sigma_{i,j} \quad (1.18)$$

the final covariance will be given by the formula:

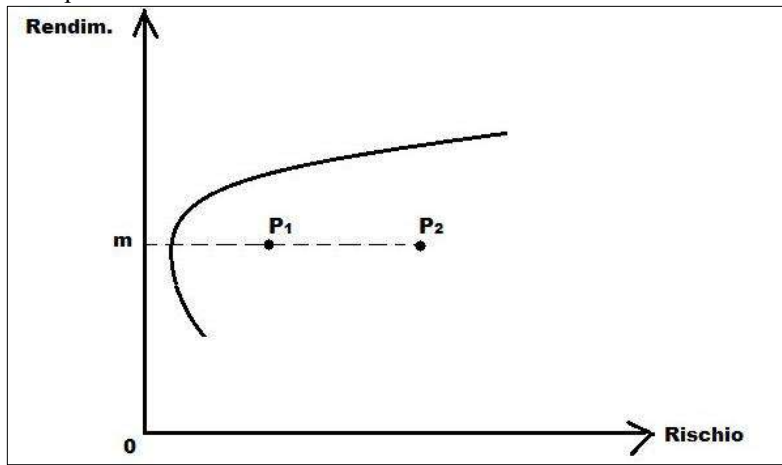
$$\text{Cov}(F_0, F(m)) = m\alpha + \beta \quad (1.19)$$

i.e. the covariance between two frontier portfolios will be the linear combination dependent on m which represents the expected return.

In the case in which short sales can be made if two securities portfolios have the same yield, it means that the covariance of both coincide, that is (Fig.4):

$$\text{Cov}(P_1, F) = \text{Cov}(P_2, F) \quad (1.20)$$

Figure 4- Mean / Variance plane and Efficient Frontier between two securities having the same performance.



Randomly choosing a portfolio like this:

$$F = x_1P_1 + x_2P_2 + x_3P_3 \quad (1.21)$$

the performance will be:

$$m = x_1m_1 + x_2m_2 + x_3m_3 \quad (1.22)$$

while the variance will be, according to the formula $\sigma_F^2 = \sum_{h,k=1}^n x_h x_k \sigma_{h,k}$, of the type:

$$\sigma_F^2 = x_1^2\sigma_1^2 + x_2^2\sigma_2^2 + x_3^2\sigma_3^2 + 2x_1x_2\sigma_{1,2} + 2x_1x_3\sigma_{1,3} + 2x_2x_3\sigma_{2,3} \quad (1.23)$$

Based on (1.20) and (1.21) then we will get:

$$\text{Cov}(P_1, F) = \text{Cov}(P_1, (x_1P_1 + x_2P_2 + x_3P_3)) = x_1\sigma_1^2 + x_2\sigma_{1,2} + x_3\sigma_{1,3} \quad (1.24)$$

$$\text{Cov}(P_2, F) = \text{Cov}(P_2, (x_1P_1 + x_2P_2 + x_3P_3)) = x_1\sigma_{1,2} + x_2\sigma_2^2 + x_3\sigma_{2,3} \quad (1.25)$$

If by absurd $\text{Cov}(P_1, F) > \text{Cov}(P_2, F)$, it means that taken another frontier portfolio of type:

$$F_k = (x_1 - k)P_1 + (x_2 - k)P_2 + x_3P_3 \quad (1.26)$$

the performance will be:

$$m_{F_k} = (x_1 - k)m_1 + (x_2 - k)m_2 + x_3m_3$$

Having assumed that P_1 and P_2 have the same performance, it means that $m_F = m_{F_k}$ while the variance of F_k will be

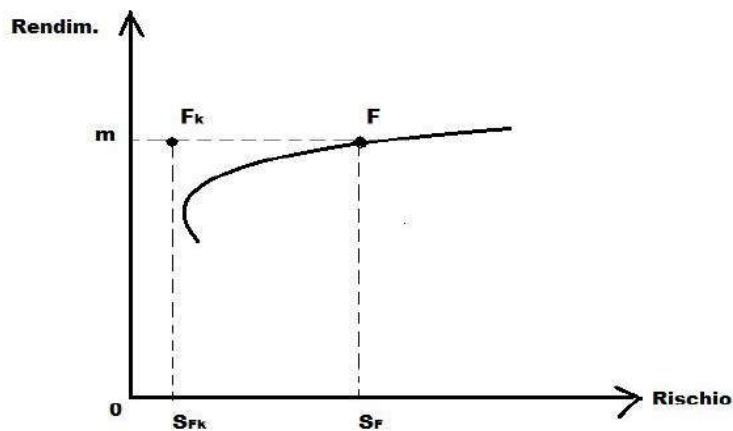
$$\begin{aligned} \sigma_{F_k}^2 = & (x_1 - k)^2 \sigma_1^2 + (x_2 - k)^2 \sigma_2^2 + x_3^2 \sigma_3^2 + 2(x_1 - k)(x_2 - k) \sigma_{1,2} + \\ & + 2(x_1 - k)x_3 \sigma_{1,3} + 2(x_2 - k)x_3 \sigma_{2,3} \end{aligned} \quad (1.27)$$

Finally we have

$$\begin{aligned} \sigma_F^2 - \sigma_{F_k}^2 = & [x_1^2 \sigma_1^2 + x_2^2 \sigma_2^2 + x_3^2 \sigma_3^2 + 2x_1x_2 \sigma_{1,2} + 2x_1x_3 \sigma_{1,3} + 2x_2x_3 \sigma_{2,3}] - \\ & - [(x_1 - k)^2 \sigma_1^2 + (x_2 + k)^2 \sigma_2^2 + x_3^2 \sigma_3^2 + 2(x_1 - k)(x_2 + k) \sigma_{1,2} + 2(x_1 - k)x_3 \sigma_{1,3} + 2(x_2 + k)x_3 \sigma_{2,3}] = \\ = & x_1^2 \sigma_1^2 + x_2^2 \sigma_2^2 + x_3^2 \sigma_3^2 + 2x_1x_2 \sigma_{1,2} + 2x_1x_3 \sigma_{1,3} + 2x_2x_3 \sigma_{2,3} - x_1^2 \sigma_1^2 - k^2 \sigma_1^2 + 2x_1k \sigma_1^2 - x_2^2 \sigma_2^2 - \\ & - k^2 \sigma_2^2 - 2x_2k \sigma_2^2 - x_3^2 \sigma_3^2 - 2x_1x_2 \sigma_{1,2} - 2x_1k \sigma_{1,2} + 2x_2k \sigma_{1,2} + 2k^2 \sigma_{1,2} - 2x_1x_3 \sigma_{1,3} + \\ & + 2x_3k \sigma_{1,3} - 2x_2x_3 \sigma_{2,3} - 2x_3k \sigma_{2,3} = \\ = & k [2(x_1 \sigma_1^2 + x_2 \sigma_{1,2} + x_3 \sigma_{1,3}) - 2(x_1 \sigma_{1,2} + x_2 \sigma_2^2 + x_3 \sigma_{2,3}) - k(\sigma_1^2 + \sigma_2^2 + 2\sigma_{1,2})] \Leftrightarrow \sigma_F^2 > \sigma_{F_k}^2 \quad (1.28) \end{aligned}$$

Graphically it can be noticed that the two portfolios F and F_k have the same performance, but different variance (Fig.5), for which F_k is not on the efficient frontier. This leads to a paradox (as expected), since it results $Cov(P_1, F) > Cov(P_2, F)$, which is impossible because from the beginning we set $Cov(P_1, F) = Cov(P_2, F)$ c.v.d.

Figure 5- Mean / Variance plane and Efficient Frontier between two security portfolios having the same performance, but different risk / variance



3. Data Envelopment Analysis applied to the European and American scenarios

For the analysis of securities portfolios in this context, the **Efficient Frontier**⁷ software was used, which implements the variance and covariance analysis algorithms developed by the Nobel Prize winner for Economics Harry Markowitz.

Through the use of these techniques and procedures, it is possible to know the variability and expected returns of any securities portfolio built on any financial asset for which volatility expectations, yield forecasts and correlation coefficients are made known. The result of the analysis will be a curve, known as the *efficient frontier*, which represents the combination of asset classes that maximize the expected return given a certain level of expected risk.

Before proceeding with the analysis of the frontier, it should be remembered that the expectations of return and risk of a portfolio are only reliable if the forecasts of return and volatility of each individual asset included in the optimization that represent the input data are reliable.

Therefore, it is also appropriate to build the efficient border on a very large number of assets, provided that there is no kind of correlation between them. In fact, the more assets there are, the lower the correlation values and the better the performance in terms of risk/return will be, thus having available a scenario in which to identify the best securities portfolio to meet the needs and personal resources of a certain investor.

In the light of the above, the objective of this work will be to compare and relate

⁷ <http://www.nonsolofondi.it/frontiera.aspx>

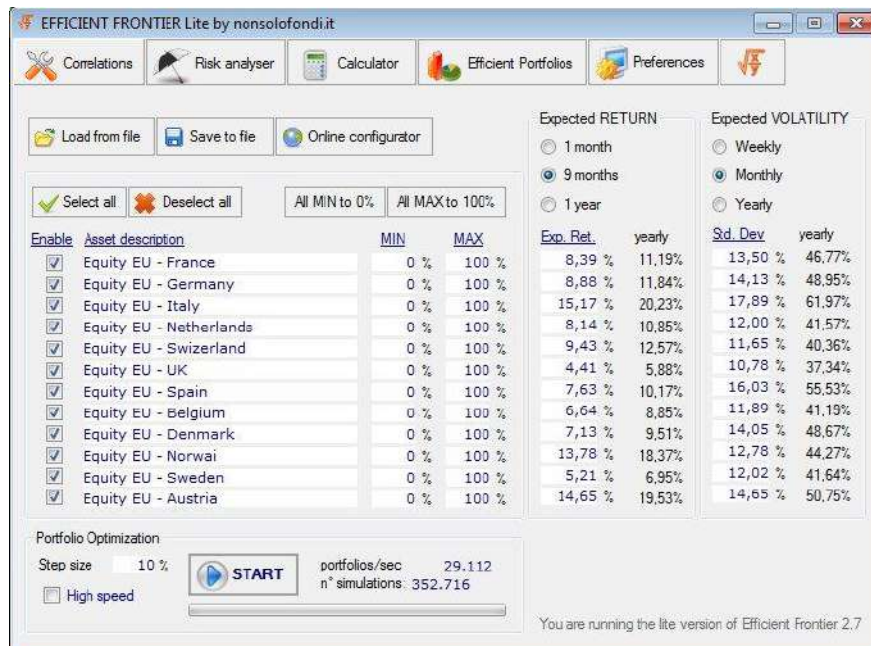
three types of market, the European and American equity with the strategic form that represents the hybrid and combined between the two. The choice of the two equity portfolios is dictated by three fundamental elements that characterize the specular peculiarities of the same.

A first feature concerns the average age of companies, so much so that in Europe in the 90s listed companies were much older than their U.S. competitors. In fact, an analysis conducted on the European market shows that the average age was 28 years, while another study conducted in America in the same period showed that the average age is reduced to ¼. This analysis has shown that the US market is a catalyst and attractor of companies and businesses that are increasingly “younger” than the European market. In addition, it was also found that existing public offers to sell securities or shares are much more frequent in Europe than in the US, where public offers to sell newly issued shares prevail.

A second characteristic concerns the role played by analysts. In the USA, companies promoting an initial public offering are subject to a so-called “quiet period” for which, in the first 40 days of issue, analysts are prohibited from issuing recommendations or research reports concerning the listed company, because the relevant information is contained in the Prospectus; subsequently, analysts play an extremely important role, as they guarantee the informative coverage of the listed company. In the European market, on the other hand, the role played by analysts is decidedly different. There are no restrictions on the quiet period: analysts can freely issue reports and analyses both during the initial public offer and in the days immediately following it.

A third characteristic concerns the major difference inherent in the concept of class action, rather rare in Europe but much more common and widespread in the United States. In fact, in the USA there are much more legal actions against individuals who behave incorrectly towards investors, while in the European market such legal actions are much more rare and unusual. For these reasons, almost all American companies tend to protect themselves through insurance cover at the expense of returns.

Figure 6 - Europe securities portfolio equity



In particular, two portfolios of securities from the European and the American area are broken down as follows (Figures 6 and 7), where we find the percentages of MIN and MAX composition of each individual asset, its expected returns and their respective volatility, from which derives the Portfolio Optimization which represents the weight of the assets in the portfolio that varies using a “increase” defined by the user (in this case it is 10%).

Figure 7 - USA securities portfolio equity

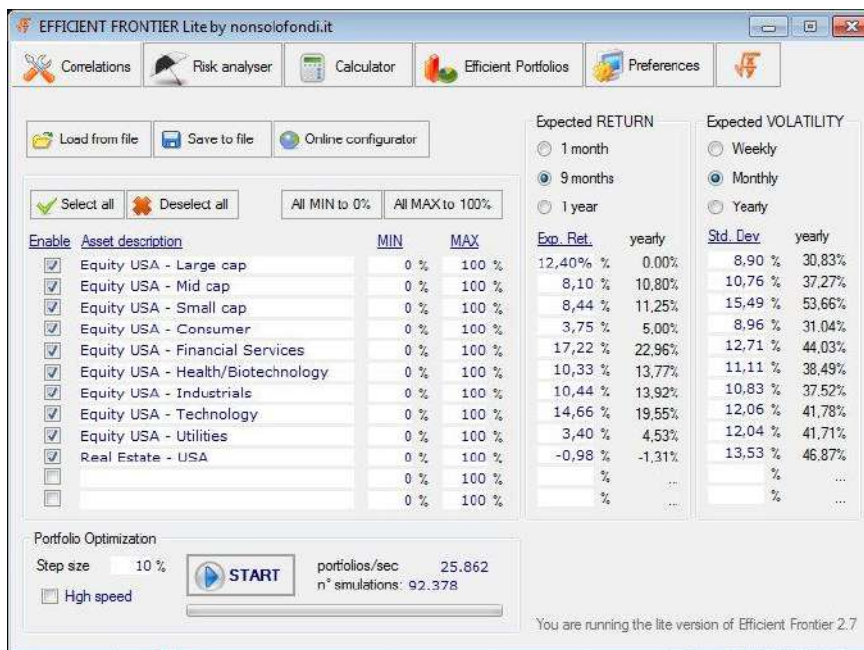
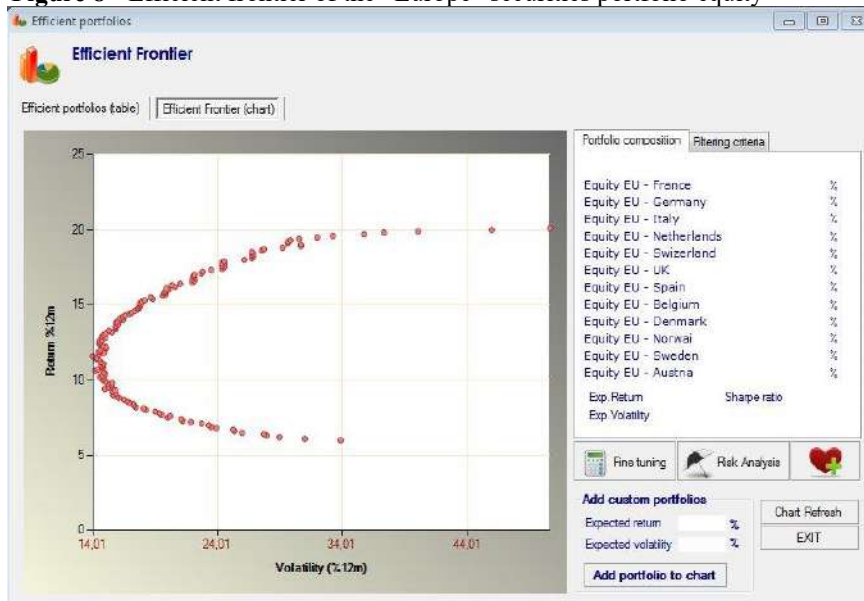
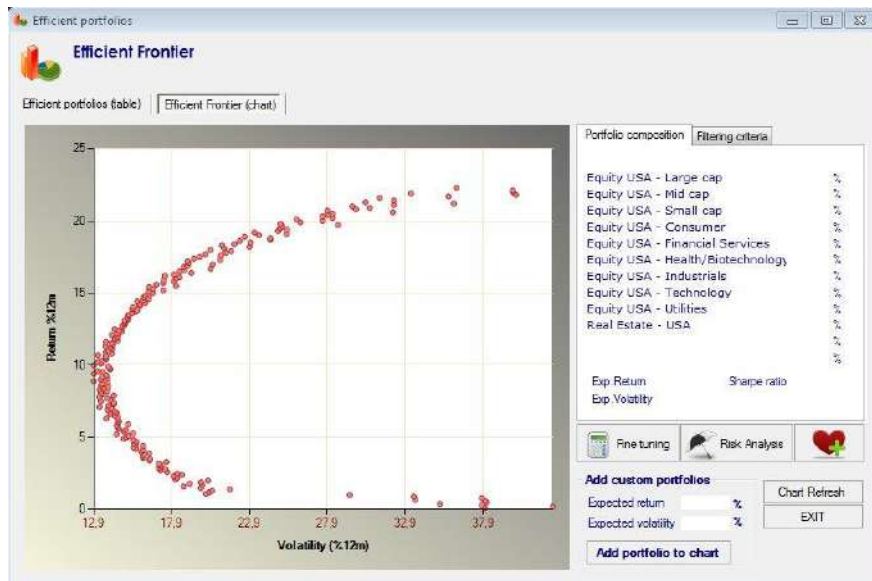


Figure 8 - Efficient frontier of the “Europe” securities portfolio equity



Consequently, the efficient borders of the European and American basins that derive from this are shown in the graphs in Figures 8 and 9.

Figure 9 - Efficient frontier of the “USA” securities portfolio equity



By default, the observation time is referred to daily prices, but this period can be extended to a week or 1 month. Beyond the month it is not possible to go, in how much it is not possible to imagine that of the formulated forecasts today they can remain valid for beyond 30 days. In any case, volatility data are always expressed on an annual basis. In addition, it is possible to act on:

- the minimum increase for the construction of the efficient border, which in this case has been set at 10%;
- the risk-free interest rate used for the calculation of the Sharpe, which in this case was considered to be 1,05%;
- the average reliability of the yield forecasts used to manage the equalization in the allocation rules set at 100%.

As can be seen, in the European case, as the risk increases, the yield does not exceed the threshold of 20%; while, in the American case, it can be seen that as volatility increases, the yield increases, reaching almost 23%.

At the end of the optimization process, for each possible level of return obtained, the ideal and optimal allocation that minimizes risk is shown in tables 1 and 2 below, illustrated as follows:

Table 1. Optimal parametric allocations that minimise the expected volatility inherent in the European area

SR	Ret	Vol	A01	A02	A03	A04	A05	A06	A07	A08	A09	A10	A11	A12
0.1458	6.00%	33.87%	0.0000	0.0000	0.0000	0.0000	0.0000	90.0000	0.0000	0.0000	0.0000	0.0000	10.0000	0.0000
0.1626	6.10%	31.01%	0.0000	0.0000	0.0000	0.0000	0.0000	80.0000	0.0000	0.0000	0.0000	0.0000	20.0000	0.0000
0.1778	6.20%	28.97%	0.0000	0.0000	0.0000	0.0000	0.0000	70.0000	0.0000	0.0000	0.0000	0.0000	30.0000	0.0000
0.1883	6.30%	27.92%	0.0000	0.0000	0.0000	0.0000	0.0000	60.0000	0.0000	0.0000	0.0000	0.0000	40.0000	0.0000
0.1925	6.40%	27.74%	0.0000	0.0000	0.0000	0.0000	0.0000	70.0000	0.0000	10.0000	0.0000	0.0000	20.0000	0.0000
0.2097	6.50%	25.98%	0.0000	0.0000	0.0000	0.0000	0.0000	60.0000	0.0000	10.0000	0.0000	0.0000	30.0000	0.0000
0.2190	6.60%	25.36%	0.0000	0.0000	0.0000	0.0000	0.0000	50.0000	0.0000	10.0000	0.0000	0.0000	40.0000	0.0000
0.2230	6.70%	25.28%	0.0000	0.0000	0.0000	0.0000	0.0000	60.0000	0.0000	20.0000	0.0000	0.0000	20.0000	0.0000
0.2401	6.80%	23.93%	0.0000	0.0000	0.0000	0.0000	0.0000	50.0000	0.0000	20.0000	0.0000	0.0000	30.0000	0.0000
0.2500	6.90%	23.50%	0.0000	0.0000	0.0000	0.0000	0.0000	50.0000	10.0000	10.0000	0.0000	0.0000	30.0000	0.0000
0.2564	7.00%	23.31%	10.0000	0.0000	0.0000	0.0000	0.0000	50.0000	0.0000	10.0000	0.0000	0.0000	30.0000	0.0000
0.2669	7.10%	22.73%	0.0000	0.0000	0.0000	0.0000	0.0000	50.0000	10.0000	20.0000	0.0000	0.0000	20.0000	0.0000
0.2824	7.20%	21.86%	0.0000	0.0000	0.0000	0.0000	0.0000	40.0000	10.0000	20.0000	0.0000	0.0000	30.0000	0.0000
0.2935	7.30%	21.23%	0.0000	0.0000	0.0000	0.0000	0.0000	40.0000	10.0000	10.0000	10.0000	0.0000	30.0000	0.0000
0.3024	7.40%	21.08%	0.0000	0.0000	0.0000	10.0000	0.0000	40.0000	10.0000	10.0000	0.0000	0.0000	30.0000	0.0000
0.3244	7.50%	20.03%	0.0000	0.0000	0.0000	10.0000	0.0000	40.0000	0.0000	20.0000	10.0000	0.0000	20.0000	0.0000
0.3266	7.60%	20.20%	0.0000	10.0000	0.0000	0.0000	0.0000	40.0000	0.0000	20.0000	10.0000	0.0000	20.0000	0.0000
0.3396	7.70%	19.53%	0.0000	0.0000	0.0000	10.0000	0.0000	40.0000	10.0000	10.0000	10.0000	0.0000	20.0000	0.0000
0.3513	7.80%	19.35%	0.0000	10.0000	0.0000	10.0000	0.0000	40.0000	0.0000	10.0000	10.0000	0.0000	20.0000	0.0000
0.3597	7.90%	19.01%	10.0000	0.0000	0.0000	10.0000	0.0000	30.0000	0.0000	10.0000	10.0000	0.0000	30.0000	0.0000
0.3788	8.00%	18.29%	0.0000	0.0000	0.0000	10.0000	10.0000	30.0000	10.0000	20.0000	10.0000	0.0000	20.0000	0.0000

Table 2. Optimal parametric allocations that minimise the expected volatility inherent in the USA area

SR	Ret	Vol	A01	A02	A03	A04	A05	A06	A07	A08	A09	A10	A11	A12
-0.0255	0.10%	38.00%	10.0000	0.0000	10.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	80.0000	0.0000	0.0000
-0.0197	0.20%	42.35%	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	90.0000	0.0000	0.0000
-0.0196	0.30%	37.81%	10.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	80.0000	0.0000	0.0000
-0.0172	0.40%	35.11%	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	30.0000	70.0000	0.0000	0.0000
-0.0136	0.50%	38.11%	0.0000	0.0000	10.0000	0.0000	0.0000	0.0000	0.0000	0.0000	10.0000	80.0000	0.0000	0.0000
-0.0124	0.60%	38.00%	0.0000	0.0000	10.0000	10.0000	0.0000	0.0000	0.0000	0.0000	0.0000	80.0000	0.0000	0.0000
-0.0101	0.70%	33.53%	10.0000	0.0000	10.0000	10.0000	0.0000	0.0000	0.0000	0.0000	0.0000	70.0000	0.0000	0.0000
-0.0054	0.80%	37.81%	0.0000	0.0000	0.0000	10.0000	0.0000	0.0000	0.0000	0.0000	0.0000	80.0000	0.0000	0.0000
-0.0036	0.90%	33.43%	10.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	10.0000	0.0000	70.0000	0.0000	0.0000
-0.0001	1.00%	29.34%	20.0000	0.0000	0.0000	0.0000	0.0000	10.0000	0.0000	0.0000	10.0000	60.0000	0.0000	0.0000
0.0047	1.10%	20.12%	50.0000	0.0000	0.0000	10.0000	0.0000	0.0000	0.0000	0.0000	20.0000	20.0000	0.0000	0.0000
0.0093	1.20%	20.30%	50.0000	0.0000	0.0000	30.0000	0.0000	0.0000	0.0000	0.0000	0.0000	20.0000	0.0000	0.0000
0.0133	1.30%	20.50%	60.0000	0.0000	0.0000	20.0000	0.0000	0.0000	0.0000	0.0000	10.0000	10.0000	0.0000	0.0000
0.0180	1.40%	21.62%	50.0000	0.0000	0.0000	0.0000	0.0000	10.0000	0.0000	0.0000	10.0000	30.0000	0.0000	0.0000
0.0235	1.50%	19.78%	30.0000	0.0000	0.0000	20.0000	0.0000	0.0000	0.0000	0.0000	20.0000	30.0000	0.0000	0.0000
0.0319	1.60%	18.66%	40.0000	0.0000	0.0000	20.0000	0.0000	0.0000	0.0000	0.0000	20.0000	20.0000	0.0000	0.0000
0.0347	1.70%	19.85%	40.0000	0.0000	0.0000	40.0000	0.0000	0.0000	0.0000	0.0000	0.0000	20.0000	0.0000	0.0000
0.0405	1.80%	19.07%	50.0000	0.0000	0.0000	30.0000	0.0000	0.0000	0.0000	0.0000	10.0000	10.0000	0.0000	0.0000
0.0449	1.90%	19.79%	40.0000	0.0000	0.0000	10.0000	0.0000	10.0000	0.0000	0.0000	10.0000	30.0000	0.0000	0.0000
0.0498	2.00%	20.23%	50.0000	0.0000	0.0000	0.0000	0.0000	0.0000	10.0000	0.0000	20.0000	20.0000	0.0000	0.0000
0.0603	2.10%	18.15%	30.0000	0.0000	0.0000	30.0000	0.0000	0.0000	0.0000	0.0000	20.0000	20.0000	0.0000	0.0000

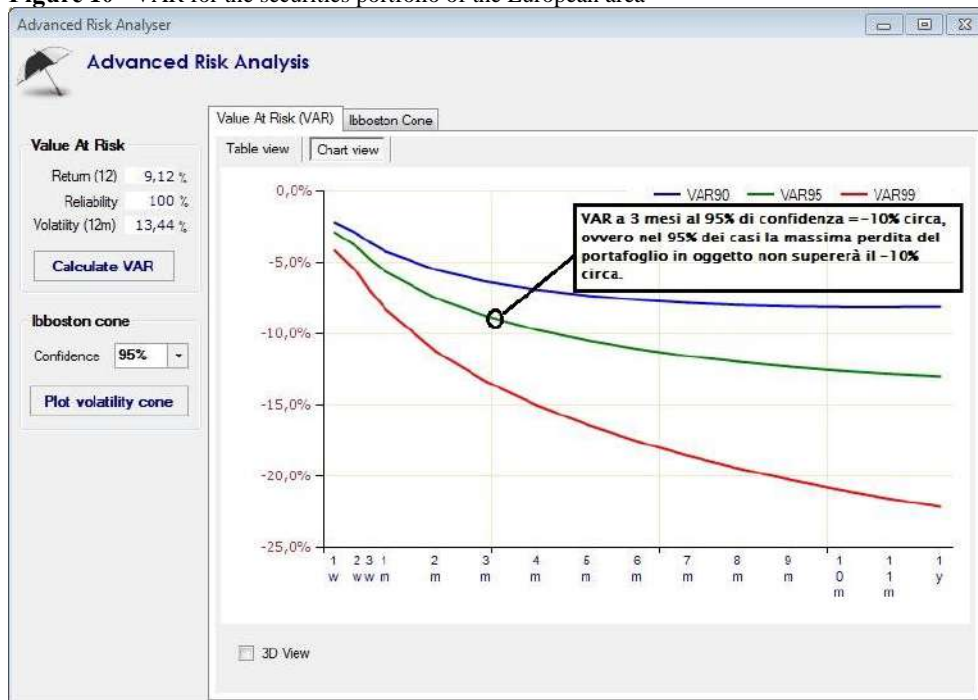
4. Analysis of the Var and its efficient border

The Value at Risk (VAR) is a statistical indicator of risk that aims to provide an indication of the extent of the range that may have the loss of a certain asset allocation in case of adverse, but not extreme, market scenarios. The VAR is therefore always associated with a probability, indicating the threshold within which any losses should remain confined, even if it is presented as a very precise number, i.e. an exact figure. In fact, although this is the best possible estimate, to be exact, the *value-at-risk* assumes that the following three conditions are met at the same time:

1. The expected return is really the average value of a hypothetical frequency distribution of returns generated by that given portfolio;
2. The volatility estimates are correct;
3. The frequency distribution of portfolio returns is a “normal” distribution.

For the European and US securities portfolios, the VAR and the volatility cone⁸ are represented by Figures 10-13.

Figure 10 - VAR for the securities portfolio of the European area



⁸ **Volatility cone:** shows the expected behavior of the portfolio over time.

Figure 11 - VAR for the securities portfolio of the USA area

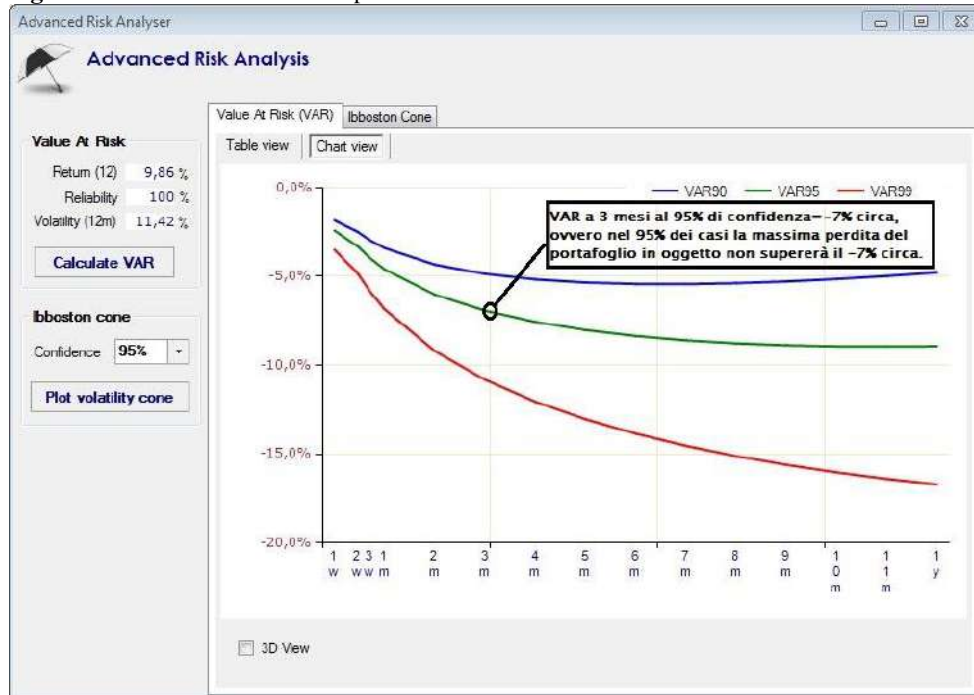


Figure 12 - Volatility Cone for the securities portfolio of the European area.

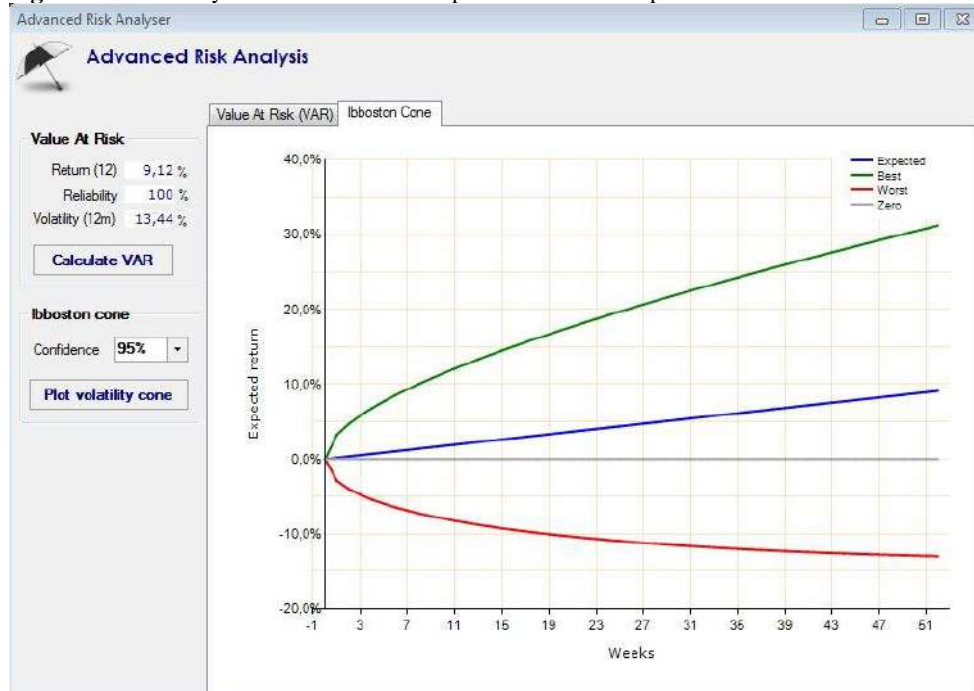
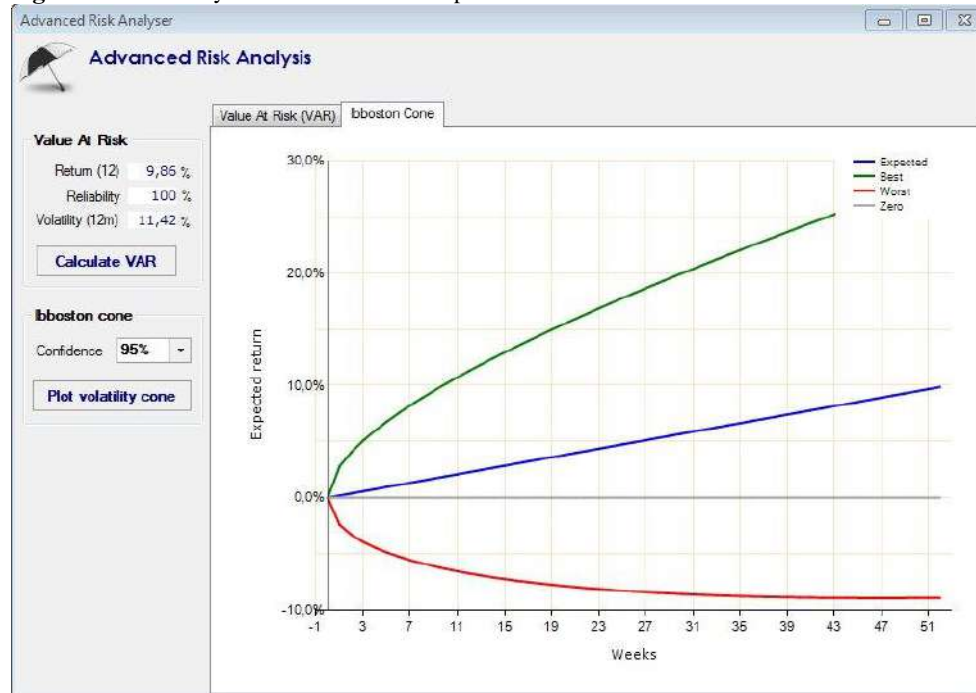


Figure 13 - Volatility Cone for the securities portfolio of the USA area.

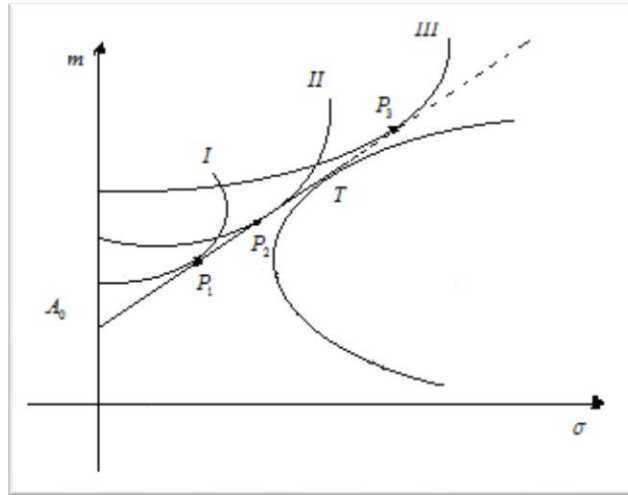
As it is possible to notice from the graphs illustrating the cones of volatility, it can be easily deduced how in both cases the greater losses will be concentrated between the third and the fifty-first week. As time passes, the volatility plays against the investor and causes that, in fact, after 3 weeks there will be a 95% probability of obtaining a return between -15% and +30% in the first case (European area), and a return between -10% and +25% in the second case (U.S. area).

As can be seen from the conicity of the two charts, it can be deduced that the greater volatility (expressed in terms of range) present in the European market is an indicator of greater risk; while in the US market, the smaller opening of the cone is a symptom of less risk. This will lead a potential risk averse investor to prefer the portfolio of securities in the American basin over the European one.

If, moreover, to this objective picture is added the subjective component (expressed by the curves of indifference that have a convexity less and less accentuated to the decrease of the aversion to the risk), and it is considered the possibility to combine the risky portfolio with one without risk, it is possible to identify the optimal portfolio on the efficient frontier given by the straight line with interception in A_0 and tangent to the function frontier of the opportunities, identifying therefore a point T on the same one that expresses the part of optimal risk. Therefore, the best portfolio

will be identified by the shifts of the indifference curves on the A_0T line, towards A_0 in the hypothesis of high risk aversion values. In the case of persons with little risk aversion, their strategy is to sell short A_0 and then finance themselves at an i_0 rate and invest in units of the risky T portfolio (Fig.14).

Figure 14 – Efficient frontier and indifference curves



Conclusions

The procedure illustrated is interesting in order to determine which is the best efficient frontier for an investor in order to maximize his premium-return in the medium/long term.

In particular, the objective is to apply the DEA to the securities of the two basins (in particular the European and the American one) that are more subject to greater transactions at the equity level. This allows the investor, who intends to capitalize in these financial instruments, to build and delineate efficient frontiers optimal and appropriate to its economic and financial possibilities. In this way, it is guaranteed, at the same time, to be able to focus on the most convenient and favourable scenario, which guarantees - at the same time - a more attractive and reasonable profit which (in the case in question) was found to be the American one, the result being the most convenient in terms of risk-expected return.

References

1. AdvisingTools: Efficient Frontier 3.0 (2018). <http://www.advisingtools.com/efficientfrontier.aspx>
2. Amirteimoori, A.: Dea efficiency analysis: Efficient and anti-efficient frontier. *Applied Mathematics and Computation* 186(1), 10-16 (2007)
3. Bisceglia, M.: Il controllo di gestione. Università degli Studi di Bari \Aldo Moro" (2015)
4. Bodnar, T., Schmid, W.: Econometrical analysis of the sample efficient frontier. *The European journal of Finance* 15(3), 317-335 (2009)
5. Broll, U., Egozcue, M., Wong, W.K., Zitikis, R.: Prospect theory, indifference curves, and hedging risks. *Applied Mathematics Research Express* 2010(2), 142-153 (2010)
6. Charnes, A., Cooper, W.W., Lewin, A.Y., Seiford, L.M.: *Data envelopment analysis: Theory, methodology, and applications*. Springer Science & Business Media (2013)
7. Charnes, A., Cooper, W.W., Rhodes, E.: Measuring the efficiency of decision making units. *European journal of operational research* 2(6), 429-444 (1978)
8. Elbannan, M.A.: The capital asset pricing model: an overview of the theory. *International Journal of Economics and Finance* 7(1), 216-228 (2015)
9. Elton, E.J., Gruber, M.J., Brown, S.J., Goetzmann, W.N.: *Modern portfolio theory and investment analysis*. John Wiley & Sons (2009)
10. Esfahani, H.N., hossein Sobhiyah, M., Yousef, V.R.: Project portfolio selection via harmony search algorithm and modern portfolio theory. *Procedia-Social and Behavioral Sciences* 226, 51-58 (2016)
11. Fabozzi, F.J., Gupta, F., Markowitz, H.M.: The legacy of modern portfolio theory. *The Journal of Investing* 11(3), 7-22 (2002)
12. Farrell, M.J.: The measurement of productive efficiency. *Journal of the Royal Statistical Society: Series A (General)* 120(3), 253-281 (1957)
13. Feldman, D., Reisman, H.: Simple construction of the efficient frontier. *European Financial Management* 9(2), 251-259 (2003)
14. Hodges, S.D., Tompkins, R.: Volatility cones and their sampling properties. *The Journal of Derivatives* 10(1), 27-42 (2002)
15. Irwin, R.J., Irwin, T.C.: A principled approach to setting optimal diagnostic thresholds: where ROC and indifference curves meet. *European Journal of Internal Medicine* 22(3), 230-234 (2011)
16. Markowitz, H.M.: Portfolio theory: as I still see it. *Annu. Rev. Financ. Econ.* 2(1), 1-23 (2010)

17. Premachandra, I., Chen, Y., Watson, J.: Dea as a tool for predicting corporate failure and success: A case of bankruptcy assessment. *Omega* 39(6), 620-626 (2011)
18. Rockafellar, R.T., Uryasev, S.: Conditional value-at-risk for general loss distributions. *Journal of banking & finance* 26(7), 1443-1471 (2002)
19. Shalit, H., Yitzhaki, S.: The mean-Gini efficient portfolio frontier. *Journal of Financial Research* 28(1), 59-75 (2005)
20. Yamai, Y., Yoshida, T., et al.: On the validity of value-at-risk: comparative analyses with expected shortfall, *Monetary and Economic Studies*, Vol. 20, No. 1, Bank of Japan, pp.57-86 (2002)

La stabilità nelle equazioni differenziali lineari

Mauro Gianfranco Bisceglia*

Dipartimento di Economia e Finanza, Università degli Studi di Bari Aldo Moro

Riassunto: Il presente lavoro non vuole avere la pretesa di affrontare il problema della stabilità nelle equazioni differenziali ordinarie di ordine k , ma di fornire una semplice trattazione del problema di una soluzione \bar{y} , appartenente all'insieme delle soluzioni $W(f, k)$ di un'equazione differenziale che risulti: stabile, asintoticamente stabile, uniformemente stabile o asintoticamente-uniformemente stabile. Pertanto, al fine di poter affrontare al meglio tale argomento, si è resa indispensabile un'introduzione su alcune questioni di base, quali la trattazione delle stesse equazioni alle differenze (ordinarie) di ordine k , con i problemi che ne derivano nell'individuazione dell'integrale generale per la risoluzione di ogni problema di Cauchy ad essa relativo. Servendosi quindi della funzione lipschitziana rispetto alla i -esima coordinata, si è giunti all'unica soluzione del problema di Cauchy. Infine si è affrontato un po' più da vicino il problema della stabilità per l'equazione differenziale lineare sia del primo che del secondo ordine a coefficienti costanti. A completamento si riporta una semplice applicazione di un modello economico.

Keywords: equazioni differenziali; equazioni alle differenze di ordine k ; soluzione stabile; asintoticamente stabile; uniformemente stabile; asintoticamente-uniformemente stabile.

1. Alcuni elementi introduttivi

Iniziamo con il considerare il seguente problema.

1.1. Siano: $k \in \mathbb{N}$, X una parte di \mathbb{R}^{k+2} ed f una funzione reale definita in X . Determinare tutte le funzioni reali u definite in un intervallo $I \subseteq \mathbb{R}$, ivi derivabili k volte, con derivata k -esima continua e tali inoltre che, $\forall x \in I$, si abbia:

$$f(x, u(x), u'(x), \dots, u^{(k)}(x)) = 0 . \quad [1.2]$$

*Autore corrispondente: maurogianfranco.bisceglia@uniba.it.

Il problema sopra formulato si denota con il simbolo:

$$f(x, y, y', \dots, y^k) = 0 \quad [1.3]$$

e si chiama *equazione alle differenze* (ordinarie) di ordine k .

Ogni funzione u soddisfacente a quanto richiesto dal problema 1.1 si chiama soluzione dell'equazione lineare [1.3] o equivalentemente integrale particolare dell'equazione [1.3].

Si dice che la [1.3] è un'equazione differenziale di ordine k del tipo normale se esiste una funzione reale g definita in R^{k+1} tale che la [1.3] possa scriversi sotto la forma:

$$y^k = g(x, y, y', \dots, y^{k-1}). \quad [1.4]$$

Si dice che la [1.4] è un'equazione differenziale lineare di ordine k se esistono $k+1$ funzioni lineari, a_1, a_2, \dots, a_k ed s , definite in un intervallo $I \subseteq R$ e tali che l'equazione [1.4] possa scriversi sotto la forma:

$$y^k = a_1(x)y^{k-1} + a_2(x)y^{k-2} + \dots + a_k(x)y + s(x), \quad [1.5]$$

le funzioni a_1, a_2, \dots, a_k risultano i coefficienti dell'equazione lineare [1.5] e la funzione s il termine noto dell'equazione lineare [1.5].

Se $\forall x \in I$ risulta $s(x) = 0$, allora l'equazione [1.5] si dice lineare omogenea di ordine k .

A volte insieme all'equazione [1.5] viene considerata anche l'equazione

$$y^k = a_1(x)y^{k-1} + a_2(x)y^{k-2} + \dots + a_k(x)y; \quad [1.6]$$

in questo caso alla [1.6] si dà il nome di equazione omogenea associata all'equazione [1.5].

Tornando all'equazione [1.3], il poterla risolvere significa trovare tutti i suoi integrali particolari.

Def. 1.7. Siano $Y \subseteq R^{k+1}$ e $Y \neq \emptyset$ e w una funzione reale definita in Y . Si dice che w è un integrale generale dell'equazione [1.3] se, per ogni elemento (c_1, c_2, \dots, c_k) di R^k per cui esiste un intervallo $J \subseteq R$ tale che $J \times \{(c_1, c_2, \dots, c_k)\} \subseteq Y$, la restituzione di $w(c_1, c_2, \dots, c_k)$ a J è un integrale particolare dell'equazione [1.3] e se, u essendo un integrale particolare dell'equazione [1.3] definito in un intervallo $I \subseteq R$, esiste un elemento (c_1, c_2, \dots, c_k) di R^k tale che $I \times \{(c_1, c_2, \dots, c_k)\} \subseteq Y$ e la restrizione di $w(c_1, c_2, \dots, c_k)$ ad I coincida con u .

Osservazione 1.8. Si noti che l'equazione differenziale [1.3] è risolta se si riesce a definire un suo integrale generale.

Relativamente all'equazione differenziale [1.4] si pongono diversi problemi, il più comune dei quali è il seguente:

Problema di Cauchy o dei valori iniziali 1.9. Sia $(x_0, y_0, y_0', \dots, y_0^{k-1})$ un elemento dell'insieme di definizione della funzione g .

Determinare, se esistono, tutti gli integrali particolari u dell'equazione [1.4], definiti in un intervallo di R a cui x_0 appartiene e tali che risulti $u(x_0) = y_0, u'(x_0) = y_0', \dots, u^{k-1}(x_0) = y_0^{k-1}$.

Il problema 1.9 qui formulato si denota con il simbolo:

$$\begin{cases} y^k = g(x, y, y', \dots, y^{k-1}) \\ y(x_0) = y_0, y'(x_0) = y_0', \dots, y^{k-1}(x_0) = y_0^{k-1} \end{cases}$$

Applicazione 1.10. Siano: I un intervallo di R avente interno non vuoto, f una funzione reale definita in I ed ivi continua, ed a un elemento di I .

Possiamo far vedere che la funzione:

$$(x, c_1, c_2, \dots, c_k) \in I \times R^k \rightarrow \sum_{h=1}^k c_{k+1-h} \frac{(x-a)^{h-1}}{(h-1)!} + \int_a^x \frac{(x-t)^{k-1}}{(k-1)!} f(t) dt \quad [1.11]$$

è un integrale generale dell'equazione differenziale:

$$y^k = f(x) \quad [1.12]$$

Dim. Per $k = 1$ quanto asserito risulta evidentemente vero.

Dimostriamo quindi l'asserto per induzione. Per questo supponiamo che la funzione [1.11] sia un integrale generale dell'equazione [1.12] e facciamo vedere che la funzione:

$$(x, c_1, c_2, \dots, c_k, c_{k+1}) \in I \times R^{k+1} \rightarrow \sum_{h=1}^{k+1} c_{k+2-h} \frac{(x-a)^{h-1}}{(h-1)!} + \int_a^x \frac{(x-t)^k}{k!} f(t) dt$$

è un integrale generale dell'equazione differenziale:

$$y^{k+1} = f(x) \quad [1.13]$$

A tale scopo osserviamo che, in virtù di quanto sopra supposto, per ogni elemento (c_1, c_2, \dots, c_k) di R^k , la funzione:

$$x \in I \rightarrow \sum_{h=1}^k c_{k+1-h} \frac{(x-a)^{h-1}}{(h-1)!} + \int_a^x \frac{(x-t)^{k-1}}{(k-1)!} f(t) dt$$

è la derivata prima di una delle soluzioni dell'equazione differenziale [1.13], quindi la funzione:

$$x \in I \rightarrow \sum_{h=1}^{k+1} c_{k+2-h} \frac{(x-a)^{h-1}}{(h-1)!} + \int_a^x \left(\int_a^y \frac{(y-t)^{k-1}}{(k-1)!} f(t) dt \right) dy$$

è, per ogni numero reale c_{k+1} , una soluzione dell'equazione differenziale [1.13].

Si osserva ancora che essendo:

$$\begin{aligned} \int_a^x \left(\int_a^y \frac{(y-t)^{k-1}}{(k-1)!} f(t) dt \right) dy &= \frac{1}{(k-1)!} \int_a^x \left(\sum_{i=0}^{k-1} \binom{k-1}{i} (-1)^i y^{k-i-1} \int_a^y t^i f(t) dt \right) dy = \\ &= \frac{1}{(k-1)!} \left\{ \left[\sum_{i=0}^{k-1} \binom{k-1}{i} \frac{(-1)^i y^{k-i}}{k-i} \int_a^y t^i f(t) dt \right]_a^x - \int_a^x \sum_{i=0}^{k-1} \binom{k-1}{i} \frac{(-1)^i y^{k-i}}{k-i} y^i f(y) dy \right\} = \\ &= \frac{1}{(k-1)!} \left\{ \sum_{i=0}^{k-1} \binom{k-1}{i} \frac{(-1)^i y^{k-i}}{k-i} \int_a^x t^i f(t) dt - \sum_{i=0}^{k-1} \binom{k-1}{i} \frac{(-1)^i}{k-i} \int_a^x t^k f(t) dt \right\}, \end{aligned}$$

e risultando $\binom{k-1}{i} = \frac{1}{k} \binom{k}{i}$, è:

$$\int_a^x \left(\int_a^y \frac{(y-t)^{k-1}}{(k-1)!} f(t) dt \right) dy = \frac{1}{k!} \left\{ \sum_{i=0}^{k-1} \binom{k}{i} (-1)^i x^{k-1} \int_a^x t^i f(t) dt - \sum_{i=0}^{k-1} \binom{k}{i} (-1)^i \int_a^x t^k f(t) dt \right\}$$

ed osservando infine che è:

$$\sum_{i=0}^{k-1} \binom{k}{i} (-1)^i + (-1)^k \binom{k}{k} = 0,$$

come volevasi, risulta:

$$\begin{aligned} \int_a^x \left(\int_a^y \frac{(y-t)^{k-1}}{(k-1)!} f(t) dt \right) dy &= \frac{1}{k!} \left\{ \sum_{i=0}^{k-1} \binom{k}{i} (-1)^i x^{k-1} \int_a^x t^i f(t) dt - (-1)^k \binom{k}{k} \int_a^x t^k f(t) dt \right\} = \\ &= \frac{1}{k!} \int_a^x \sum_{i=0}^k \binom{k}{i} (-1)^i x^{k-1} t^i f(t) dt = \int_a^x \frac{(x-t)^k}{k!} f(t) dt. \end{aligned}$$

Osservazione 1.14. Si noti che se di un'equazione differenziale si conosce un integrale generale, allora è possibile risolvere ogni problema di Cauchy ad essa relativo.

2. Un Teorema di esistenza ed unicità per la soluzione del problema di Cauchy

Iniziamo con il porre la seguente definizione:

Def. 2.1. Siano I_1, I_2, \dots, I_k k intervalli non vuoti di R , f una funzione reale definita in $I_1 \times I_2 \times \dots \times I_k$, ed $i \in \{1, 2, \dots, k\}$.

Si dice che f è lipschitziana rispetto alla i -esima coordinata se esiste un numero reale positivo L tale che, per ogni coppia $((x_1, x_2, \dots, x'_i, \dots, x_k), (x_1, x_2, \dots, x''_i, \dots, x_k))$ di elementi di $I_1 \times I_2 \times \dots \times I_k$ le cui coordinate diverse dalla i -esima sono uguali, risulta:

$$|f(x_1, x_2, \dots, x'_i, \dots, x_k) - f(x_1, x_2, \dots, x''_i, \dots, x_k)| \leq L |x'_i - x''_i|.$$

Il numero reale positivo L si chiama il coefficiente di Lipschitz di f rispetto alla i -esima coordinata.

Di immediata verifica è la seguente proposizione:

2.2. Siano: I_1, I_2, \dots, I_k k intervalli non vuoti di R , f una funzione reale definita in $I_1 \times I_2 \times \dots \times I_k$ ed $i \in \{1, 2, \dots, k\}$.

Se f è parzialmente derivabile rispetto alla coordinata I_1, I_2, \dots, I_k e se la derivata parziale f_{x_i} è limitata, allora f è lipschitziana rispetto alla i -esima coordinata.

Dim. Infatti se $(x_1, x_2, \dots, x'_i, \dots, x_k)$ ed $(x_1, x_2, \dots, x''_i, \dots, x_k)$ sono due elementi di $I_1 \times I_2 \times \dots \times I_k$ aventi coordinate diverse dalla i -esima uguali, in virtù del teorema di Lagrange, esiste un elemento y di $[\min\{x'_i, x''_i\}, \max\{x'_i, x''_i\}]$ tale che risulti:

$$f(x_1, x_2, \dots, x'_i, \dots, x_k) - f(x_1, x_2, \dots, x''_i, \dots, x_k) = f_{x_i}(x_1, x_2, \dots, y, \dots, x_k)(x'_i - x''_i);$$

detto quindi L l'estremo superiore di $|f_{x_i}|$ in $I_1 \times I_2 \times \dots \times I_k$, è immediato rendersi conto che, come volevasi, risulta:

$$|f(x_1, x_2, \dots, x'_i, \dots, x_k) - f(x_1, x_2, \dots, x''_i, \dots, x_k)| \leq L|x'_i - x''_i|.$$

Dimostriamo ora la seguente proposizione:

2.3. Siano: I un intervallo chiuso, limitato e non vuoto di R , f una funzione reale definita in $I \times R$ ed ivi continua, $x_0 \in I$ ed $y_0 \in R$.

Se f è lipschitziana rispetto alla seconda coordinata, allora esiste una ed una sola soluzione definita in I del problema di Cauchy:

$$\begin{cases} y' = f(x, y) \\ y(x_0) = y_0 \end{cases} \quad [2.4]$$

Dim. Sia $u_1 : x \in I \rightarrow y_0 + \int_{x_0}^x f(t, y_0) dt$ e $\forall n \in N - \{1\}$; poniamo:

$$u_n : x \in I \rightarrow y_0 + \int_{x_0}^x f(t, u_{n-1}(t)) dt$$

Cominciamo con l'osservare che se denotiamo con L il coefficiente di Lipschitz di f rispetto alla seconda coordinata e poniamo $M = \max_{x \in I} |u_1(x) - y_0|$, allora,

$\forall n \in N$ e $\forall x \in I$, risulta:

$$|u_{n+1}(x) - u_n(x)| \leq M \frac{(L|x - x_0|)^n}{n!} \quad [2.5]$$

Infatti essendo:

$$|u_2(x) - u_1(x)| = \left| \int_{x_0}^x f(t, u_1(t)) - f(t, y_0) dt \right| \leq \left| \int_{x_0}^x |u_{n+1}(t) - u_n(t)| dt \right| \leq ML|x - x_0|,$$

la [2.5] è vera per $n = 1$; se supponiamo la [2.5] vera per n , allora $\forall x \in I$ risulta:

$$\begin{aligned} |u_{n+2}(x) - u_{n+1}(x)| &= \left| \int_{x_0}^x f(t, u_{n+1}(t)) - f(t, u_n(t)) dt \right| \leq \left| \int_{x_0}^x L|u_{n+1}(t) - u_n(t)| dt \right| \leq \\ &\leq \frac{LML^n}{n!} \left| \int_{x_0}^x |t - x_0|^n dt \right| = M \frac{(L|x - x_0|)^{n+1}}{(n+1)!} \end{aligned}$$

e quindi la [2.5] resta dimostrata per il principio di induzione.

Dalla [2.5] consegue che, se diciamo d la dimensione dell'intervallo I , $\forall x \in I$ e $\forall n \in N$ risulta:

$$|u_{n+1}(x) - u_n(x)| \leq M \frac{(Ld)^n}{n!}. \quad [2.6]$$

Se n ed h sono due numeri interi positivi è facile dimostrare che, $\forall x \in I$, risulta:

$$|u_{n+h}(x) - u_n(x)| = \sum_{i=1}^h |u_{n+i}(x) - u_{n+i-1}(x)| \leq \sum_{i=1}^h |u_{n+i}(x) - u_{n+i-1}(x)|.$$

Conseguentemente, in virtù della [2.6], è:

$$|u_{n+h}(x) - u_n(x)| \leq \sum_{i=1}^h M \frac{(Ld)^{n+1-i}}{(n+1-i)!} = \sum_{i=1}^{n+h-1} M \frac{(Ld)^i}{i!} - \sum_{i=1}^{n-1} M \frac{(Ld)^i}{i!}$$

ed essendo la serie dei termini generale $M \frac{(Ld)^n}{n!}$ convergente, in virtù del criterio di

convergenza di Cauchy, possiamo affermare che per ogni numero reale positivo ε esiste un numero $m \in N$ tale che $\forall n \in N$, con $n > m$, $\forall h \in N$ e $\forall x \in I$ risulta:

$$|u_{n+h}(x) - u_n(x)| \leq \varepsilon, \quad [2.7]$$

e questo, sempre in virtù del criterio di convergenza di Cauchy, è quanto dire che $\forall x \in I$ la successione $(u_n(x))_{n \in N}$ converge; posto quindi $u(x) = \lim_n u_n(x)$, dalla

[2.7] consegue che, per ogni numero reale positivo ε , $\exists m \in N$ tale che $\forall x \in I$ e $\forall n \in N$, con $n > m$, risulta:

$$|u(x) - u_n(x)| \leq \varepsilon. \quad [2.8]$$

Facciamo ora osservare che $\forall x \in I$ è:

$$\lim_n \int_{x_0}^x f(t, u_{n-1}(t)) dt = \int_{x_0}^x f(t, u(t)) dt \quad [2.9]$$

e osserviamo anche che $\forall x \in I$ si ha:

$$\left| \int_{x_0}^x f(t, u(t)) dt - \int_{x_0}^x f(t, u_{n-1}(t)) dt \right| \leq \int_{x_0}^x |f(t, u(t)) - f(t, u_{n-1}(t))| dt \leq \int_{x_0}^x L |u(t) - u_{n-1}(t)| dt;$$

quindi, in virtù della [2.8], possiamo asserire che per ogni numero reale positivo ε , $\exists m \in N$ tale che $\forall x \in I$ e $\forall n \in N$, con $n > m$, è:

$$\left| \int_{x_0}^x f(t, u(t)) dt - \int_{x_0}^x f(t, u_{n-1}(t)) dt \right| \leq Ld\varepsilon$$

dimostrando la [2.9], ed essendo $u_n(x) = y_0 + \int_{x_0}^x f(t, u_{n-1}(t)) dt \quad \forall n \in N$ e $\forall x \in I$,

consegue che $\forall x \in I$ è:

$$u(x) = y_0 + \int_{x_0}^x f(t, u(t)) dt; \quad [2.10]$$

da questa uguaglianza consegue che $\forall x \in I$ è $u'(x) = u(x, u(x))$: quest'ultima uguaglianza e la [2.10] ci assicurano che la funzione u è una soluzione del problema [2.4].

Infine dimostriamo che u è l'unica soluzione del problema [2.4].

Per questo supponiamo che esista un'altra soluzione v del problema [2.4] e chiamiamola $\bar{x} \in I$, ove risulti $u(\bar{x}) \neq v(\bar{x})$; supponendo $\bar{x} > x_0$, se si pone $\bar{\bar{x}} = \sup\{x \in [x_0, \bar{x}] : u(\bar{x}) = v(\bar{x})\}$ è evidente che risulta: $x_0 \leq \bar{\bar{x}} \leq \bar{x}$, $u(\bar{\bar{x}}) = v(\bar{\bar{x}})$ e $|u(x) - v(x)| > 0$, $\forall x \in]\bar{\bar{x}}, \bar{x}]$, e quindi essendo:

$$u(x) = u(\bar{\bar{x}}) + \int_{\bar{\bar{x}}}^x f(t, u(t)) dt \quad \text{e} \quad v(x) = v(\bar{\bar{x}}) + \int_{\bar{\bar{x}}}^x f(t, v(t)) dt$$

è

$$|u(x) - v(x)| = \left| \int_{\bar{\bar{x}}}^x (f(t, u(t)) - f(t, v(t))) dt \right| \leq \int_{\bar{\bar{x}}}^x L |u(t) - v(t)| dt;$$

conseguentemente $\forall x \in I$ è:

$$|u(x) - v(x)| \leq \int_{\bar{\bar{x}}}^x |u(t) - v(t)| dt \quad [2.11]$$

Sia ora $\delta \in R_+$ minore di $\bar{x} - \bar{\bar{x}}$ e tale inoltre che risulti che $L\delta < 1$, se diciamo quindi H il massimo della funzione $|u - v|$ in $[\bar{\bar{x}}, \bar{\bar{x}} + \delta]$ dalla [2.11] consegue

che $\forall x \in [\bar{x}, \bar{x} + \delta]$ è $|u(x) - v(x)| \leq L\delta H < H$, e ciò è assurdo: quindi u è, come volevasi, l'unica soluzione del problema [2.4].

Dalla proposizione 2.3 consegue la seguente altra proposizione:

2.12. Siano: J un intervallo non vuoto di R , f una funzione reale definita in $J \times R$ ed ivi continua, $x_0 \in J$ e $y_0 \in R$.

Se per ogni intervallo chiuso, limitato e non vuoto $I \subseteq J$ la restrizione della funzione f a $I \times R$ è lipschitziana rispetto alla seconda coordinata, allora esiste una ed una sola soluzione del problema di Cauchy [2.4] definita in J .

Dim. Siano: $(I_n)_{n \in N}$ una successione crescente di intervalli di R chiusi, limitati e non vuoti per cui risulti: $J = \bigcup_{n \in N} I_n$ e $m \in N$ tale che $x_0 \in I$ in virtù della proposizione 2.3, $\forall n \in N$ esiste una ed una sola soluzione u_n del problema di Cauchy [2.4] definita in I_{m+n} .

Osserviamo ora che, sempre in virtù della proposizione 2.3, se $\{r, s\} \in N$ e se $r > s$ allora la restrizione di u_r ad I_{m+s} è uguale ad u_s , posto quindi:

$$\bar{u}_n : x \in J \rightarrow \begin{cases} u_n(x) & \text{se } x \in I_{m+n} \\ 0 & \text{se } x \in J - I_{m+n} \end{cases}$$

è immediato rendersi conto che, $\forall x \in J$ l'insieme $\{w \in R, \exists n \in N : w = \bar{u}_n(x)\}$ è costituito da un numero finito di elementi, conseguentemente la successione $(\bar{u}_n(x))_{n \in N}$ è convergente.

È ora facile convincersi che la funzione:

$$u : x \in J \rightarrow \lim_n \bar{u}_n(x) \in R$$

è l'unica soluzione del problema di Cauchy [2.4] definita in J .

Allo scopo di dare una generalizzazione della proposizione 2.3 e quindi pure della proposizione 2.12 cominciamo con il formulare la seguente applicazione:

Applicazione 2.13. Siano: $k \in N$, $X \subseteq R^{k+2}$ con $X \neq \emptyset$ ed f_1, f_2, \dots, f_k k funzioni reali definite in X .

Determinare tutte le funzioni vettoriali (u_1, u_2, \dots, u_k) definite in un intervallo $I \subseteq R$ aventi ogni coordinata derivabile con derivata continua e tali inoltre che, $\forall x \in I$ e $\forall i \in N$ con $i \leq k$ si abbia:

$$u'_i(x) = f_i(x, u_1(x), u_2(x), \dots, u_k(x)).$$

L'applicazione sopra formulata si denota col simbolo:

Siamo ora in grado di dimostrare la seguente proposizione:

2.19. Siano: J un intervallo non vuoto di R , f una funzione reale definita in $J \times R^k$ ed ivi continua, $x_0 \in J$ ed $(y_0, y'_0, \dots, y_0^{k-1}) \in R^k$.

Se per ogni intervallo chiuso, limitato e non vuoto $I \subseteq J$ la restrizione della funzione f ad $I \times R^k$ è lipschitziana rispetto a tutte le coordinate, tranne al più la prima, allora esiste una ed una sola soluzione definita in J del problema di Cauchy

$$\begin{cases} y^k = f(x, y, y', \dots, y^{k-1}) \\ y(x_0) = y_0, y'(x_0) = y'_0, \dots, y^{k-1}(x_0) = y_0^{k-1} . \end{cases}$$

Dim. L'asserto è conseguenza immediata della proposizione 2.18 quando si osserva che le soluzioni del problema di Cauchy considerato nell'enunciato della proposizione 2.19 sono le prime coordinate delle soluzioni del seguente altro problema di Cauchy:

$$\begin{cases} y'_1 = y_2 \\ y'_2 = y_3 \\ \dots\dots\dots \\ y'_{k-1} = y_k \\ y'_k = f(x, y_1, y_2, \dots, y_k) \\ y_1(x_0) = y_0, y_2(x_0) = y'_0, \dots, y_k(x_0) = y_0^{k-1} \end{cases} .$$

3. Sulla Stabilità nelle equazioni differenziali

Prima di dare qualche cenno sull'importanza del problema della stabilità è utile fare qualche premessa.

Cominciamo col porre la seguente definizione:

Def. 3.1. Siano $h \in N$, $S \subseteq R^h$ con $S \neq \emptyset$, $T \neq \emptyset$, x_0 un elemento in cui possa effettuarsi il limite su S , g una funzione reale definita in $S \times T$ e $V \in \hat{R}$.

Si dice che il limite di $g(x, y)$ al tendere di x ad x_0 è uguale a V uniformemente rispetto a y in T e si scrive:

$$(uni.y \in T) \lim_{x \rightarrow x_0} g(x, y) = V$$

Se $\forall I_V, \exists J_{x_0}$ tale che $\forall (x, y) \in ((S - \{x_0\}) \cap J_{x_0}) \times T$ $g(x, y) \in I_V$.

In questo paragrafo $k \in N$, $t \in \{-\infty\} \cup R$ ed f una funzione reale definita in $]t, +\infty[\times R$ soddisfacente le ipotesi del teorema 2.19.

$\forall x_0 \in]t, +\infty[$ è $u(x) = w(x, c_1(x_0, y_0, y'_0, \dots, y_0^{k-1}), c_2(x_0, y_0, y'_0, \dots, y_0^{k-1}), \dots, c_k(x_0, y_0, y'_0, \dots, y_0^{k-1}))$, conseguentemente la funzione $\phi: (x, x_0, y_0, \dots, y_0^{k-1}) \in]t, +\infty[^2 \times R^k \rightarrow w(x, c_1(x_0, y_0, y'_0, \dots, y_0^{k-1}), c_2(x_0, y_0, y'_0, \dots, y_0^{k-1}), \dots, c_k(x_0, y_0, y'_0, \dots, y_0^{k-1}))$ come è facile verificare, non dipende dall'integrale generale w considerato; quindi essa, a differenza dell'integrale generale dell'equazione differenziale [3.3], è univocamente determinata, inoltre ϕ è derivabile parzialmente k volte rispetto alla prima coordinata x risultando, $\forall i \in \{1, 2, \dots, k\}$ e $\forall x \in]t, +\infty[$:

$$\frac{\partial^i \phi(x, x_0, y_0, y'_0, \dots, y_0^{k-1})}{\partial x^i} = D^i u(x)$$

È ora possibile dimostrare che:

3.5. ϕ e le sue prime $k-1$ derivate parziali rispetto alla prima coordinata x sono k funzioni continue in $]t, +\infty[^2 \times R^k$.

Per una più comoda e spedita trattazione conviene porre la seguente definizione:

Def. 3.6. Siano: u una funzione reale $k-1$ derivabile, definita in $]t, +\infty[$ e $x_0 \in]t, +\infty[$.

Con il simbolo u_0 denotiamo l'elemento di $R^k : (u(x_0), u'(x_0), \dots, u^{k-1}(x_0))$.

Siamo ora in grado di porre la seguente definizione:

Def. 3.7. Sia $\bar{y} \in W(f, k)$.

Si dice che \bar{y} è stabile per l'equazione [3.3] se, $\forall x \in]t, +\infty[$ e $\forall i \in \{1, 2, \dots, k-1\}$, risulta:

$$\begin{aligned} & (\text{uni. } x \in]x_0, +\infty[) \lim_{y_0 \rightarrow \bar{y}_0} |\phi(x, x_0, y_0) - \phi(x, x_0, \bar{y}_0)| = 0 \\ & (\text{uni. } x \in]x_0, +\infty[) \lim_{y_0 \rightarrow \bar{y}_0} \left| \frac{\partial^i \phi(x, x_0, y_0)}{\partial x^i} - \frac{\partial^i \phi(x, x_0, \bar{y}_0)}{\partial x^i} \right| = 0 \end{aligned}$$

Si dice che \bar{y} è asintoticamente stabile per l'equazione [3.3] se, essendo stabile, $\forall x_0 \in]t, +\infty[\exists I_{\bar{y}_0}$ tale che $\forall i \in \{1, 2, \dots, k-1\}$ risulta:

$$\begin{aligned} & (\text{uni. } y_0 \in I_{\bar{y}_0}) \lim_{x \rightarrow +\infty} |\phi(x, x_0, y_0) - \phi(x, x_0, \bar{y}_0)| = 0 \\ & (\text{uni. } y_0 \in I_{\bar{y}_0}) \lim_{x \rightarrow +\infty} \left| \frac{\partial^i \phi(x, x_0, y_0)}{\partial x^i} - \frac{\partial^i \phi(x, x_0, \bar{y}_0)}{\partial x^i} \right| = 0 \end{aligned}$$

Si dice che \bar{y} è uniformemente stabile per l'equazione [3.3] se $\forall i \in \{1, 2, \dots, k-1\}$ risulta:

d'altra parte se denotiamo con A il denominatore della matrice dei coefficienti del sistema [3.10] e, $\forall (i, j) \in \{1, 2, \dots, k\}^2$ denotiamo con $A_{i,j}$ dell'elemento di posto (i, j) nella matrice dei coefficienti del sistema [3.10] e poniamo $\bar{a}_{i,j} = \frac{A_{i,j}}{A}$ è subito visto che $\forall i \in \{1, 2, \dots, k\}$ è:

$$c_1^0 = (\bar{y}(x_0) - v(x_0))\bar{a}_{1,i} + (\bar{y}'(x_0) - v'(x_0))\bar{a}_{2,i} + \dots + (\bar{y}^{k-1}(x_0) - v^{k-1}(x_0))\bar{a}_{k,i}$$

conseguentemente $\forall x \in]t, +\infty[$ è:

$$\begin{aligned} \bar{y}(x) = & (\bar{y}(x_0) - v(x_0)) \sum_{i=1}^k \bar{a}_{1,i} u_i(x) + (\bar{y}'(x_0) - v'(x_0)) \sum_{i=1}^k \bar{a}_{2,i} u_i(x) + \dots \\ & + (\bar{y}^{k-1}(x_0) - v^{k-1}(x_0)) \sum_{i=1}^k \bar{a}_{k,i} u_i(x) \end{aligned}$$

quindi se y è un'altra soluzione dell'equazione differenziale [3.9], $\forall (x, x_0) \in]t, +\infty[$ risulta:

$$\begin{aligned} \phi(x, x_0, y_0) - \phi(x, x_0, \bar{y}_0) = & (y(x_0) - \bar{y}(x_0)) \sum_{i=1}^k \bar{a}_{1,i} u_i(x) + (y'(x_0) - \bar{y}'(x_0)) \\ & \sum_{i=1}^k \bar{a}_{2,i} u_i(x) + \dots + (y^{k-1}(x_0) - \bar{y}^{k-1}(x_0)) \sum_{i=1}^k \bar{a}_{k,i} u_i(x) = \phi(x, x_0, y_0 - \bar{y}_0) \quad [3.11] \end{aligned}$$

e di qui, quando si osservi che la funzione $x \in]t, +\infty[\rightarrow \phi(x, x_0, y_0 - \bar{y}_0)$ è la soluzione del seguente problema di Cauchy:

$$\begin{cases} z^k = a_1(x)z^{k-1} + a_2(x)z^{k-2} + \dots + a_k(x)z \\ z(x_0) = y(x_0) - \bar{y}(x_0), z'(x_0) = y'(x_0) - \bar{y}'(x_0), \dots, z^{k-1}(x_0) = y^{k-1}(x_0) - \bar{y}^{k-1}(x_0) \end{cases}$$

e che $\forall (x, x_0) \in]t, +\infty[$ se u è la funzione identicamente nulla in $]t, +\infty[$, è $\phi(x, x_0, u_0) = 0$ consegue quanto asserito.

Dall'uguaglianza [3.11] consegue la seguente altra proposizione:

3.12. Condizione necessaria e sufficiente affinché la funzione identicamente nulla in $]t, +\infty[$ asintoticamente-uniformemente stabile per l'equazione omogenea associata all'equazione [3.9] è che essa sia uniformemente stabile ed asintoticamente stabile per l'equazione omogenea associata all'equazione [3.9].

Osserviamo ora un po' più da vicino il problema della stabilità per l'equazione differenziale lineare del primo ordine. Per questo se a è una funzione reale, definita e continua in $]t, +\infty[$ è facile far vedere che se y è una soluzione dell'equazione

$$y' = a(x)y \tag{3.13}$$

allora $\forall (x, x_0) \in]t, +\infty[^2$ è:

$$\begin{aligned} \phi(x, x_0, y_0) - \phi(x, x_0, \bar{y}_0) &= (y(x_0) - \bar{y}(x_0)) \sum_{i=1}^k \bar{a}_{1,i} u_i(x) + (y'(x_0) - \bar{y}'(x_0)) \\ \phi(x, x_0, y_0) &= y(x_0) e^{A(x) - A(x_0)} \end{aligned} \tag{3.14}$$

essendo A una primitiva della funzione a .

Cominciamo con il dimostrare la seguente proposizione:

3.15. Condizione necessaria e sufficiente affinché la funzione identicamente nulla in $]t, +\infty[$ sia stabile per l'equazione [3.13] è che $\exists \bar{x} \in]t, +\infty[$ ed $\exists h \in \mathbb{R}$, tali che, $\forall x \in]\bar{x}, +\infty[$, risulta $A(x) - A(\bar{x}) = \int_{[\bar{x}, x]} a dm_1 \leq h$.

Dim. Facciamo dapprima vedere che la condizione è sufficiente. Per questo notiamo che esistono un elemento $\bar{x} \in]t, +\infty[$ ed un $h \in \mathbb{R}$ tali che risulti: $A(x) - A(\bar{x}) \leq h \quad \forall x \in]\bar{x}, +\infty[$, se $x_0 \in]t, +\infty[$, risulta:

$$\begin{aligned} \int_{[x_0, x]} a dm_1 &= \int_{[\bar{x}, x]} a dm_1 - \int_{[\bar{x}, x_0]} a dm_1 \leq h - \int_{[\bar{x}, x_0]} a dm_1 \quad \text{se } \bar{x} \leq x_0 \leq x \\ &= \int_{[x_0, x]} a dm_1 \leq \int_{[x_0, \bar{x}]} |a| dm_1 \quad \text{se } x_0 \leq x \leq \bar{x} \\ \int_{[x_0, x]} a dm_1 &= \int_{[x_0, \bar{x}]} a dm_1 + \int_{[\bar{x}, x]} a dm_1 \leq \int_{[\bar{x}, x_0]} a dm_1 + h \quad \text{se } x_0 \leq \bar{x} \leq x \end{aligned}$$

quindi se poniamo $I = [\min\{\bar{x}, x_0\}, \max\{\bar{x}, x_0\}]$ ed $h_0 = |h| + \int_{[I]} |a| dm_1$,

$\forall x \in [x_0, +\infty[$ è $A(x) - A(x_0) < h$, conseguentemente, $\forall x \in [x_0, +\infty[$ e per ogni soluzione y dell'equazione [3.13], è $|\phi(x, x_0, y_0)| \leq |y(x_0)| e^{h_0}$ e ciò dimostra quanto volevasi.

Dimostriamo ora che la condizione è necessaria. Per questo supponiamo che pur essendo la funzione identicamente nulla in $]t, +\infty[$ stabile per l'equazione [3.13] $\forall (\bar{x}, h) \in]t, +\infty[\times \mathbb{R} \quad \exists x \in [\bar{x}, +\infty[$ tale che risulti $A(x) - A(\bar{x}) > h$, in virtù di queste ipotesi, detto $x_0 \in]t, +\infty[$, $\forall d \in \mathbb{R}_+$ diciamo $x \in [x_0, +\infty[$ tale che sia $A(x) - A(\bar{x}) > \log \frac{2}{d}$ e diciamo D la soluzione dell'equazione [3.13] che in x_0

assume valore $\frac{d}{2}$, è immediato rendersi conto che risulta $|\phi(x, x_0, D_0)| > 1$, quindi abbiamo in fine dimostrato che $\exists x_0 \in]t, +\infty[$ ed $\exists \varepsilon = 1$, tale che $\forall d \in R_+$ $\exists x \in [x_0, +\infty[$ ed una soluzione D dell'equazione [3.13] tale che, pur essendo nulla in $]t, +\infty[$ non è stabile per l'equazione [3.13].

L'asserto resta così completamente dimostrato.

Dalla proposizione 3.15 ora dimostrata consegue la seguente altra:

3.16. Condizione necessaria e sufficiente affinché la funzione identicamente nulla in $]t, +\infty[$ sia asintoticamente stabile per l'equazione [3.13] è che sia stabile per l'equazione [3.13] e risulti: $\lim_{x \rightarrow +\infty} A(x) = -\infty$.

Dim. Se la funzione identicamente nulla in $]t, +\infty[$ è stabile l'equazione [3.13] essendo:

$$\lim_{x \rightarrow +\infty} |\phi(x, x_0, y_0)| = |y(x_0)| \lim_{x \rightarrow +\infty} e^{A(x) - A(x_0)}$$

se è $\lim_{x \rightarrow +\infty} A(x) = -\infty$ allora, se $\forall \delta \in R_+$, risulta:

$$(\text{univ}, y(x_0) \in]-\delta, \delta]) \lim_{x \rightarrow +\infty} |\phi(x, x_0, y_0)| = 0$$

e ciò dimostra la sufficienza della condizione.

Facciamo ora vedere che la condizione è necessaria. Per questo supponiamo che pur essendo la funzione identicamente nulla in $]t, +\infty[$ asintoticamente stabile per l'equazione [3.13] non sia stabile per l'equazione [3.13] o non risulti $\lim_{x \rightarrow +\infty} A(x) = -\infty$. In ognuna di queste eventualità non è difficile dimostrare che, contro l'ipotesi, la funzione identicamente nulla in $]t, +\infty[$ non è asintoticamente stabile per l'equazione [3.13].

L'asserto resta così completamente dimostrato.

È ora immediato dimostrare che:

3.17. Condizione necessaria e sufficiente affinché la funzione identicamente nulla in $]t, +\infty[$ sia uniformemente stabile per l'equazione [3.13] è che $\exists H \in R$ tale che, $\forall x_0 \in]t, +\infty[$ e $\forall x \in [x_0, +\infty[$, risulti: $A(x) - A(x_0) \leq H$.

Dalla proposizione 3.17, ricordando la proposizione 3.12, consegue che:

3.18. Condizione necessaria e sufficiente affinché la funzione identicamente nulla in $]t, +\infty[$ sia asintoticamente-uniformemente stabile per l'equazione [3.13] è che uniformemente stabile per l'equazione [3.13] e risulti: $\lim_{x \rightarrow +\infty} A(x) = -\infty$.

Si noti che in particolare si ha:

3.19. Sia a una funzione costante in $]t, +\infty[$.

Condizione necessaria e sufficiente affinché la funzione identicamente nulla in $]t, +\infty[$ sia uniformemente stabile per l'equazione [3.13] è che risulti: $a \leq 0$.

Condizione necessaria e sufficiente affinché la funzione identicamente nulla in $]t, +\infty[$ sia asintoticamente-uniformemente stabile per l'equazione [3.13] è che risulti: $a < 0$.

Qui di seguito studieremo un po' più da vicino il problema della stabilità per l'equazione differenziale lineare del secondo ordine a coefficienti costanti. Per questo siano $\{a_1, a_2\} \in R$, è facile far vedere che se y è una soluzione dell'equazione

$$y'' = a_1 y' + a_2 y \quad [3.20]$$

allora $\forall (x, x_0) \in]t, +\infty[$ è:

$$\phi(x, x_0, y_0) = y(x_0) \frac{\alpha e^{\beta(x-x_0)} - \beta e^{\alpha(x-x_0)}}{\alpha - \beta} + y'(x_0) \frac{e^{\alpha(x-x_0)} - e^{\beta(x-x_0)}}{\alpha - \beta} \quad [3.21']$$

$$\phi_x(x, x_0, y_0) = \alpha \beta y(x_0) \frac{e^{\beta(x-x_0)} - e^{\alpha(x-x_0)}}{\alpha - \beta} + y'(x_0) \frac{\alpha e^{\alpha(x-x_0)} - \beta e^{\beta(x-x_0)}}{\alpha - \beta} \quad [3.21'']$$

se è $\Delta = a_1^2 + 4a_2 > 0$ e si è posto che $\alpha = \frac{a_1 - \sqrt{\Delta}}{2}$ ed $\beta = \frac{a_1 + \sqrt{\Delta}}{2}$

$$\phi(x, x_0, y_0) = (y(x_0)(1 - \alpha(x-x_0)) + y'(x_0)(x-x_0)) e^{\alpha(x-x_0)} \quad [3.22']$$

$$\phi_x(x, x_0, y_0) = (-a^2 y(x_0)(x-x_0) + y'(x_0)(1 + \alpha(x-x_0))) e^{\alpha(x-x_0)} \quad [3.22'']$$

se è $\Delta = a_1^2 + 4a_2 = 0$ e si è posto che $\alpha = \frac{a_1}{2}$ ed $\beta = \frac{a_1 + \sqrt{\Delta}}{2}$

$$\begin{aligned} \phi(x, x_0, y_0) &= \\ &= (y(x_0)(\beta \cos \beta(x-x_0) - \alpha \operatorname{sen} \beta(x-x_0)) + y'(x_0) \operatorname{sen} \beta(x-x_0)) \frac{e^{\alpha(x-x_0)}}{\beta} \end{aligned} \quad [3.23']$$

$$\begin{aligned} \phi_x(x, x_0, y_0) &= \\ &= (-y(x_0)(\alpha^2 + \beta^2) \operatorname{sen} \beta(x-x_0) + y'(x_0)(\beta \cos \beta(x-x_0) + \alpha \operatorname{sen} \beta(x-x_0))) \frac{e^{\alpha(x-x_0)}}{\beta} \end{aligned} \quad [3.23'']$$

se è $\Delta = a_1^2 + 4a_2 < 0$ e si è posto che $\alpha = \frac{a_1}{2}$ ed $\beta = \frac{\sqrt{\Delta}}{2}$.

Osservando ora che sussiste la seguente proposizione:

3.24. Siano α e β i numeri reali che figurano nelle uguaglianze [3.21'], [3.21''], [3.22'], [3.22''], [3.23'] ed [3.23''].

Condizione necessaria e sufficiente affinché sia $\alpha \leq 0$ e $\beta \leq 0$ è che risulti $a_1 \leq 0$ e $a_2 \leq 0$.

Condizione necessaria e sufficiente affinché sia $\alpha < 0$ e $\beta < 0$ è che risulti $a_1 < 0$ e $a_2 < 0$.

Possiamo enunciare il seguente teorema:

3.25. Condizione necessaria e sufficiente affinché la funzione identicamente nulla in $]t, +\infty[$ sia uniformemente stabile per l'equazione [3.20] è che risulti: $a_1 \leq 0$ e $a_2 \leq 0$ e $a_1 + a_2 < 0$.

Condizione necessaria e sufficiente affinché la funzione identicamente nulla in $]t, +\infty[$ sia asintoticamente-uniformemente stabile per l'equazione [3.20] è che risulti: $a_1 < 0$ e $a_2 < 0$.

Passiamo ora allo studio della stabilità per l'equazione di Eulero-Lagrange.

A tale proposito osserviamo che se $(a, b, p, q) \in]0, +\infty[\times R^3$ e si prende un $x \in]-\frac{b}{a}, +\infty[$, un integrale generale dell'equazione differenziale

$$(ax+b)^2 y'' = p(ax+b)y' + qy \quad [3.26]$$

è

- la funzione

$$(x, c_1, c_2) \in]-\frac{b}{a}, +\infty[\times R^2 \rightarrow c_1(ax+b)^{\frac{\alpha}{a}} + c_2(ax+b)^{\frac{\beta}{a}},$$

se è $(p+a)^2 + 4p > 0$ ed α e β sono le due radici reali e distinte dell'equazione:

$$t^2 - (p+a)t - q = 0; \quad [3.27]$$

- la funzione:

$$(x, c_1, c_2) \in]-\frac{b}{a}, +\infty[\times R^2 \rightarrow \left(c_1 + c_2 \frac{1}{a} \log(ax+b) \right) (ax+b)^{\frac{\alpha}{a}}$$

se è $(p+a)^2 + 4p = 0$ ed α è la radice reali dell'equazione [3.27];

- la funzione:

$$(x, c_1, c_2) \in]-\frac{b}{a}, +\infty[\times R^2 \rightarrow \left(c_1 \cos\left(\frac{\beta}{a} \log(ax+b)\right) + c_2 \operatorname{sen}\left(\frac{\beta}{a} \log(ax+b)\right) \right) (ax+b)^{\frac{\alpha}{a}}$$

se è $(p+a)^2 + 4p < 0$ ed $\alpha + i\beta$ e $\alpha - i\beta$ sono due radici complesse coniugate dell'equazione [3.27].

Dalla proposizione 3.25 consegue facilmente la seguente altra:

3.28. Condizione necessaria e sufficiente affinché la funzione identicamente nulla in $\left] -\frac{b}{a}, +\infty \right[$ sia uniformemente stabile per l'equazione [3.26] è che risulti: $p+a \leq 0$, $q \leq 0$ e $p+a+q < 0$.

Condizione necessaria e sufficiente affinché la funzione identicamente nulla in $\left] -\frac{b}{a}, +\infty \right[$ sia asintoticamente-uniformemente stabile per l'equazione [3.26] è che risulti: $p+a < 0$ e $q < 0$.

4. Una semplice applicazione ad un problema economico

Il modello moltiplicatore- acceleratore di Samuelson.

In questo modello si suppone che il consumo C calcolato in un dato periodo sia proporzionale al reddito nazionale Y calcolato nel periodo precedente, quindi $\forall n \in N$, è:

$$C_{n+1} = bY_n \quad [4.1]$$

dove b è una costante positiva minore di uno, detta *propensione marginale al consumo*.

L'investimento I ha due componenti: l'investimento indotto I' che calcolato in un periodo è proporzionale all'incremento che il consumo subisce passando dal periodo precedente al periodo preso in esame, al periodo considerato e l'investimento autonomo I'' che si suppone costante ed il suo valore lo denoteremo con h , quindi, $\forall n \in N$, è:

$$I_{n+1} = a(C_{n+1} - C_n) + h \quad [4.2]$$

dove a è una costante positiva detta coefficiente di accelerazione.

In fine la condizione di equilibrio in questo modello è espressa, $\forall n \in N$ dall'uguaglianza:

$$Y_n = C_n + I_n \quad [4.3]$$

Dalle [4.3], [4.2] e [4.1] si ha:

$$Y_{n+2} = b(a+1)Y_{n+1} - abY_n + h \quad [4.4]$$

quindi il problema del calcolo di Y_n è ricondotto al calcolo del valore in n della soluzione di un qualunque problema di Cauchy relativo all'equazione [4.4], per esempio se sono noti i valori del reddito nazionale nei primi due periodi, diciamoli rispettivamente Y_1^0 ed Y_2^0 , basta risolvere il seguente problema dei valori iniziali:

$$\begin{cases} Y_{n+2} = b(a+1)Y_{n+1} - abY_n + h \\ Y_1 = Y_1^0 \text{ ed } Y_2 = Y_2^0 \end{cases}$$

per cui una soluzione generale dell'equazione [4.4] è la funzione:

$$(c_1, c_2) \in \mathbb{R}^2 \rightarrow c_1 (\xi^n)_{n \in \mathbb{N}} + c_2 (\zeta^n)_{n \in \mathbb{N}} + \left(\frac{h}{1-b} c_N(n) \right)_{n \in \mathbb{N}} \quad [4.5]$$

se è

$$b > \frac{4a}{(a+1)^2}, \quad \xi = \frac{b(a+1) - \sqrt{b^2(a+1)^2 - 4ab}}{2}$$

e

$$\zeta = \frac{b(a+1) + \sqrt{b^2(a+1)^2 - 4ab}}{2};$$

è la funzione:

$$(c_1, c_2) \in \mathbb{R}^2 \rightarrow c_1 \left(\left(\frac{2a}{a+1} \right)^n \right)_{n \in \mathbb{N}} + c_2 \left(n \left(\frac{2a}{a+1} \right)^n \right)_{n \in \mathbb{N}} + \left(\frac{h}{1-b} c_N(n) \right)_{n \in \mathbb{N}} \quad [4.6]$$

se è $b = \frac{4a}{(a+1)^2}$, ed è la funzione:

$$(c_1, c_2) \in \mathbb{R}^2 \rightarrow c_1 \left((\sqrt{ab})^n \cos n\alpha \right)_{n \in \mathbb{N}} + c_2 \left((\sqrt{ab})^n \sin n\alpha \right)_{n \in \mathbb{N}} + \left(\frac{h}{1-b} c_N(n) \right)_{n \in \mathbb{N}} \quad [4.7]$$

se è: $b < \frac{4a}{(a+1)^2}$ ed $\alpha \in]-\pi, \pi[$ tale che:

$$\cos \alpha = \frac{b(a+1)}{2\sqrt{ab}} \quad \text{e} \quad \sin \alpha = \frac{\sqrt{4ab - b^2(a+1)^2}}{2\sqrt{ab}},$$

conseguentemente $\forall n \in \mathbb{N}$, è:

$$Y_n = \frac{1}{\zeta - \xi} \left(\left(\zeta \left(Y_1^0 - \frac{h}{1-b} \right) - \left(Y_2^0 - \frac{h}{1-b} \right) \right) \zeta^{n-1} + \left(\left(Y_2^0 - \frac{h}{1-b} \right) - \xi \left(Y_1^0 - \frac{h}{1-b} \right) \right) \zeta^{n-1} \right) + \frac{h}{1-b}$$

se è $b > \frac{4a}{(a+1)^2}$, è:

$$Y_n = \left(\frac{2a(2-n)}{a+1} Y_1^0 + (n-1) Y_2^0 + \frac{h(a(n-3) - n + 1)}{(1-b)(a+1)} \right) \left(\frac{2a}{a+1} \right)^{n-2} + \frac{h}{1-b}$$

se è $b = \frac{4a}{(a+1)^2}$, è:

$$Y_n = (\sqrt{ab})^{n-2} \left(\left(2\sqrt{ab} Y_1^0 \cos \alpha - Y_2^0 + \frac{h(1 - 2\sqrt{ab} \cos \alpha)}{(1-b)} \right) \cos n\alpha + \left(Y_2^0 \cos \alpha - \sqrt{ab} Y_1^0 \cos 2\alpha + \frac{h(\sqrt{ab} \cos 2\alpha - \cos \alpha)}{(1-b)} \right) \frac{\text{senn} \alpha}{\text{sen} \alpha} \right) + \frac{h}{1-b}$$

se è $b < \frac{4a}{(a+1)^2}$.

Osserviamo ora che, in virtù della proposizione 3.16, la soluzione identicamente nulla è asintoticamente stabile per l'equazione $Y_{n+2} = b(a+1)Y_{n+1} - abY_n + h$ se e solo se sono soddisfatte le seguenti tre condizioni: $b < 1$, $1 + 2ab + b > 0$ ed $ab < 1$, di queste le prime due, ricordando le ipotesi in cui ci siamo posti, sono sempre soddisfatte, quindi possiamo asserire che: Condizione necessaria e sufficiente affinché la soluzione identicamente nulla sia asintoticamente stabile per l'equazione $Y_{n+2} = b(a+1)Y_{n+1} - abY_n + h$ è che risulti $ab < 1$.

Se si osserva il grafico di seguito, dove sono state riportate le curve di equazione $b = 1$, $b = \frac{1}{a}$ (la curva tratteggiata) e $b = \frac{4a}{(a+1)^2}$ (la curva continua), la fascia $]0, +\infty[\times]0, 1[$ resta divisa in quattro parti, diciamole A, B, C e D.

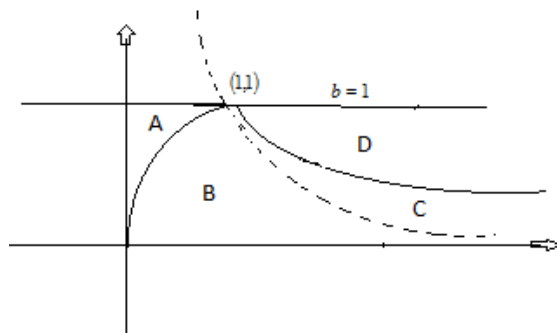


Fig. 1

Questa partizione della fascia $]0, +\infty[\times]0, 1[$ ci permette di fare la seguente analisi delle soluzioni all'equazione [4.4]:

1°) Se il punto (a, b) è interno alla zona A, allora essendo $b > \frac{4a}{(a+1)^2}$ una solu-

zione generale dell'equazione [4.4] è la funzione [4.5], inoltre, essendo anche $ab < 1$ se u è una soluzione particolare dell'equazione [4.4] è: $\lim_n u_n = \frac{h}{1-b}$.

Si osservi che essendo $0 < \xi < \zeta < 1$, se c_1 e c_2 sono due numeri reali ed u è la soluzione particolare dell'equazione [4.4] valore della funzione [4.5] in (c_1, c_2) , allora se è $c_1 \geq 0$ e $c_2 \geq 0$ risulta $\lim_n u_n = \inf_{n \in N} u_n$ e se è $c_1 \leq 0$ e $c_2 \leq 0$ risulta $\lim_n u_n = \sup_{n \in N} u_n$, se, essendo $c_1 c_2 < 0$ poniamo:

$$N' = N \cap \left] \frac{1}{\log \zeta - \log \xi} \log \left(-\frac{c_1 \log \xi}{c_2 \log \zeta} \right), +\infty \right[\quad [4.8]$$

allora se è $c_1 < 0$ e $c_2 > 0$ risulta $\lim_n u_n = \inf_{n \in N'} u_n$ e se è $c_1 > 0$ e $c_2 < 0$ risulta $\lim_n u_n = \sup_{n \in N'} u_n$.

2°) Se $b = \frac{4a}{(a+1)^2}$ e $0 < a < 1$, una soluzione generale dell'equazione [4.4] è la

funzione [4.6], inoltre essendo anche $ab < 1$ se u è una soluzione particolare dell'equazione [4.4] è: $\lim_n u_n = \frac{h}{1-b}$.

Si osservi che essendo $0 < \frac{2a}{a+1} < 1$ se c_1 e c_2 sono due numeri reali ed u è la soluzione particolare dell'equazione [4.4] valore della funzione [4.6] in (c_1, c_2) , ponendo:

$$N^0 = N \cap \left] \frac{1}{\log(a+1) - \log 2a}, +\infty \right[$$

allora se è $c_1 \geq 0$ e $c_2 \geq 0$ risulta $\lim_n u_n = \inf_{n \in N^0} u_n$ e se è $c_1 \leq 0$ e $c_2 \leq 0$ risulta $\lim_n u_n = \sup_{n \in N^0} u_n$ e ponendo:

$$N' = N \cap \left] \frac{c_2 - c_1(\log(a+1) - \log 2a)}{c_2(\log(a+1) - \log 2a)}, +\infty \right[\quad [4.9]$$

allora se è $c_1 < 0$ e $c_2 > 0$ risulta $\lim_n u_n = \inf_{n \in N'} u_n$ e se è $c_1 > 0$ e $c_2 < 0$ risulta $\lim_n u_n = \sup_{n \in N'} u_n$.

3°) Se il punto (a, b) è interno alla zona B, allora essendo $b < \frac{4a}{(a+1)^2}$ una soluzione generale dell'equazione [4.4] è la funzione [4.7], inoltre, essendo anche $ab < 1$, se u è una soluzione particolare dell'equazione [4.4] è: $\lim_n u_n = \frac{h}{1-b}$ con i valori di u che si avvicinano ad $\frac{h}{1-b}$ con un movimento oscillatorio smorzato.

Si osservi che se c_1 e c_2 sono due numeri reali non entrambi nulli ed u è la soluzione particolare dell'equazione [4.4] valore della funzione [4.7] in (c_1, c_2) , se diciamo $\gamma \in]-\pi, \pi]$ tale che $\cos \gamma = \frac{c_1}{\sqrt{c_1^2 + c_2^2}}$ e

$\sin \gamma = \frac{c_2}{\sqrt{c_1^2 + c_2^2}}$, $\forall n \in N$ è:

$$u_n = (\sqrt{ab})^n \sqrt{c_1^2 + c_2^2} \cos(n\alpha - \gamma) + \frac{h}{1-b}$$

conseguentemente i valori di u avvicinandosi ad $\frac{h}{1-b}$ oscillano e $\forall n \in N$,

è:

$$u_n \in \left[\frac{h}{1-b} - (\sqrt{ab})^n \sqrt{c_1^2 + c_2^2}, \frac{h}{1-b} + (\sqrt{ab})^n \sqrt{c_1^2 + c_2^2} \right]$$

4°) Se $a > 1$ e $b = \frac{1}{a}$ allora una soluzione generale dell'equazione [4.4] è la funzione [4.7], inoltre, se c_1 e c_2 sono due numeri reali non entrambi nulli ed u è la soluzione particolare dell'equazione [4.4] valore della funzione [4.7] in (c_1, c_2) , possiamo asserire che $\exists \lim_n u_n$ ed i valori di u oscillano intorno ad

$\frac{h}{1-b}$ con oscillazioni non costanti però tali che, $\forall n \in N$ si abbia:

$$u_n \in \left[\frac{h}{1-b} - |c_1| - |c_2|, \frac{h}{1-b} - \|c_1\| - \|c_2\| \right] \cup \left[\frac{h}{1-b} + \|c_1\| - \|c_2\|, \frac{h}{1-b} + |c_1| + |c_2| \right]$$

5°) Se il punto (a, b) è interno alla zona C, allora essendo $b < \frac{4a}{(a+1)^2}$ una soluzione generale dell'equazione [4.4] è la funzione [4.7], inoltre, essendo anche $ab > 1$, possiamo asserire che se c_1 e c_2 sono due numeri reali non entrambi nulli ed u è la soluzione particolare dell'equazione [4.4] valore della funzione [4.7] in (c_1, c_2) $\exists \lim_n u_n$ ed ivi i valori di u oscillano intorno ad $\frac{h}{1-b}$ con oscillazioni esaltate.

6°) Se è $a > 1$ e $b = \frac{4a}{(a+1)^2}$ allora una soluzione generale dell'equazione [4.4] è la funzione [4.6], inoltre, essendo $\frac{2a}{a+1} > 1$, e strettamente crescente la successione $\left(n \left(\frac{2a}{a+1} \right)^n \right)_{n \in \mathbb{N}}$, possiamo asserire che se c_1 e c_2 sono due numeri reali non entrambi nulli ed u è la soluzione particolare dell'equazione [4.4] valore della funzione [4.6] in (c_1, c_2) , risulta:

$$\lim_n u_n = \begin{cases} c_1(+\infty) & \text{se è } c_2 = 0 \\ c_2(+\infty) & \text{se è } c_2 \neq 0 \end{cases}.$$

Si osservi che se è $c_1 \geq 0$ e $c_2 \geq 0$ risulta $\lim_n u_n = \sup_{n \in \mathbb{N}} u_n$ e se è $c_1 \leq 0$ e $c_2 \leq 0$ risulta $\lim_n u_n = \inf_{n \in \mathbb{N}} u_n$ e se N' denota l'insieme definito dalla [4.9], se è $c_2 > 0$ risulta $\lim_n u_n = \sup_{n \in N'} u_n$ e se è $c_2 < 0$ risulta $\lim_n u_n = \inf_{n \in N'} u_n$.

7°) Se il punto (a, b) è interno alla zona D, allora essendo $b > \frac{4a}{(a+1)^2}$ una soluzione generale dell'equazione [4.4] è la funzione [4.5], inoltre, essendo $1 < \xi < \zeta$, se c_1 e c_2 sono due numeri reali non entrambi nulli ed u è la soluzione particolare dell'equazione [4.4] valore della funzione [4.5] in (c_1, c_2) risulta:

$$\lim_n u_n = \begin{cases} c_1(+\infty) & \text{se è } c_2 = 0 \\ c_2(+\infty) & \text{se è } c_2 \neq 0 \end{cases}$$

Si osservi che se è $c_1 \geq 0$ e $c_2 \geq 0$ risulta $\lim_n u_n = \sup_{n \in N} u_n$ e se è $c_1 \leq 0$ e $c_2 \leq 0$ risulta $\lim_n u_n = \inf_{n \in N} u_n$ e se N' denota l'insieme definito dalla [4.8], se è $c_2 > 0$ risulta $\lim_n u_n = \sup_{n \in N'} u_n$ e se è $c_2 < 0$ risulta $\lim_n u_n = \inf_{n \in N'} u_n$.

Riferimenti bibliografici

- Albano L. (1998). *Lezioni di Matematica Generale*, Cacucci Editore, Bari.
- Binmore K., Davies J. (2004). *Calcolo Differenziale di più variabili*, Casa Editrice Ambrosiana, Milano, ISBN 978-88-408-1295-4.
- De Giuli, M. E.; Giorgi, G.; Maggi, M.; Magnani, U. (2008). *Matematica per l'economia e la finanza*. Zanichelli, Bologna.
- Fedele N. (2002). *Corso di Analisi Matematica, Integrazione Equazioni Differenziali*, Volume II- Parte II, Liguori Editore, Napoli, ISBN 978-88-207-3473-2.
- Guerraggio A. (2004). *Matematica*, Mondadori, Milano, ISBN 978-88-424-9614-6.
- Maddalena L. (2009). *Matematica*, Giappichelli Editore, Torino, ISBN 978-88-348-9624-2.
- Ritelli D., Bergamini M., Trifone A. (2005). *Fondamenti di Matematica*, Zanichelli, Bologna, ISBN 978-88-08-21910-0.
- Sansone G., Conti R. (1956). *Equazioni differenziali non lineari*, Cremonese, Roma.



UNIVERSITÀ
DEGLI STUDI DI BARI
ALDO MORO

DIPARTIMENTO DI
ECONOMIA E FINANZA

PDF finito di comporre
il 29 dicembre 2020

ISBN 978-88-6629-023-0



9 788866 290230