**MICHELANGELO MISURACA | GERMANA SCEPI | MARIA SPANO**

# PROCEEDINGS OF THE
# 16TH INTERNATIONAL CONFERENCE
# ON STATISTICAL ANALYSIS OF TEXTUAL DATA
## VOLUME 1



**VADISTAT**
SIMONA BALBI

# Analysis of the public debate on DDL Zan on Twitter: an application of the Structural Topic Model

4 authors:

Maria Gabriella Grassia
University of Naples Federico II
46 PUBLICATIONS   335 CITATIONS

SEE PROFILE

Marina Marino
University of Naples Federico II
61 PUBLICATIONS   424 CITATIONS

SEE PROFILE

Rocco Mazza
University of Naples Federico II
22 PUBLICATIONS   17 CITATIONS

SEE PROFILE

Agostino Stavolo
University of Naples Federico II
15 PUBLICATIONS   7 CITATIONS

SEE PROFILE

# Analysis of the public debate on DDL Zan on Twitter: an application of the Structural Topic Model

Maria Gabriella Grassia, Marina Marino,
Rocco Mazza, Agostino Stavolo

University of Naples Federico II – mariagabriella.grassia@unina.it
marina.marino@unina.it rocco.mazza@unina.it
agostino.stavolo@unina.it

**Abstract**
Social media play a key role in analysing the public debate on political issues. The amount of information available and the ability to share opinions on various topics has led to an increase in the production of user-generated content. Indeed, data from social media platforms such as Twitter offer new possibilities because it provides a public arena to understand feedback and opinions on certain political topics. In this article, we propose an analysis of the debate on Twitter regarding the DDL Zan, the bill that aims to toughen the penalties of crimes and discrimination against homosexuals, transsexuals, women, and disabled people. We define a strategy to extract the main topics of the public debate considering the role assumed by the so-called influencers within the Twitter community. For this reason, we chose to use the Structural Topic Model (STM) including in the model the covariates describing the user's profiles of interest. This allows to identify the most relevant topics and to assess the impact of covariates in the model to study how covariates impact on them.

**Keywords:** Structural Topic Model; DDL Zan; Topic Modeling

## 1. Introduction

Social media plays a key role in analyzing public debate on policy issues, because user-generated contents (Chatziadam et al., 2020) offers a great opportunity to automatically detect opinions and trends on topics of interest. In fact, microblogging platforms such as Twitter are a great way for users to share ideas, opinions, and comments about products, services. Detecting topics in online content facilitates the identification of emerging social trends and the analysis of public reactions to policies (Sobkowicz et al., 2012). Therefore, we propose a strategy to extract the main topics of public debate by considering the role assumed by so-called influencers within the Twitter community. We present a

case study based on the use of the Structural Topic Model (STM) on a debated political event: the "DDL Zan". The bill includes measures to prevent and combat discrimination and violence on the grounds of sex, gender, sexual orientation, gender identity, and disability. The possibility of analysing public discussion on Twitter through a particular type of topic modeling allows to assess the impact of covariates in the model to study how covariates may influence the words representing a topic.

## 2. Literature Review

Twitter functions as a public sphere because everyone can see tweets, hashtags allow for organization around topics, and retweets allow for rapid dissemination (Marichal et al.,2020). On the platform, the dissemination of information is controlled through influencers (Colleoni et al., 2014). The possibility of sharing opinions and comments through tweets and the ability to interact on the platform becomes an important area of analysis for researchers. Castellò et. al (2021) performed a quantitative analysis of message interactivity as a measure of interactivity: endorsement and amplification. In this context, topic modeling facilitates understanding the conversations between people in online communities, as well as extracting useful patterns and understanding from their interactions in addition to what they share on social media websites (Hamed et al., 2019). It is used for analyzing textual content that uncovers the latent thematic structure in document collections for identifying emerging topics. In the literature, the focus of using topic modeling in recent years has been on the production of short texts on social media. Modeling in this domain is often more challenging due to character limits (e.g., tweets consisting of 280 characters). Structural topic modeling (STM) is a variant of LDA that is designed to precisely address this limitation. For this reason, STM was applied in numerous jobs to analyse online discourse (Kwon et al. 2019; He et al.,2020).

## 3. Structural Topic Model

The structural topic model (STM) introduced by Roberts et al. (2014) allows us to extract latent topics within our corpus by including a set of covariates associated with the document collection of interest. STM encompasses document metadata as covariates into the prior distributions for document-topic proportions and topic-word distributions, including additional information. The model can be decomposed into three sub-models: (a) *a topical prevalence model*, which controls how documents are allocated to topics as a function of covariates; (b) *a topical content model*, which controls words frequency in each topic as a function of covariates; (c) *a core language model*, which combines the two sources of variation to produce the actual words in each document. Formally, documents are indexed as d ∈ {1....D}, while words are indexed as n ∈ {1…N}. The

observations $w_{d;n}$ are occurrences of words from a vocabulary indexed by $v \in$ {1…V}. The number of topics is set by the researcher and indexed by $k \in$ {1… K}. The document-level additional information is represented by two matrices. X is the matrix with topical prevalence covariates, while Y is the matrix with topical content covariates. The rows of these matrices - each representing a vector of covariates for a given document - are denoted by $x_d$ and $y_d$, respectively. The process starts by drawing the document-level attention to each topic from a logistic-normal generalised linear model, based on a vector of covariates $\mathbf{x_d}$ and considering a *P x (K – 1)* matrix of coefficients *r* for the topical prevalence and a *(K-1) X (K-1)* $_d$, the topic-specific distribution over words is formed by representing each topic *k* with the V - dimensional baseline word distribution **m** (Airoldi et al. 2004), the topic specific deviation $K_k^{(t)}$, the covariate group deviation $K_{yd}^{(c)}$ and the interaction between the two just described. Finally, for each word in a document, topic assignment based on the document-specific distribution over topics and word assignment to a topic is drawn from multinomial models. The latter represents the core language model.

## 4. Data

We extracted status from Twitter using the open academy access API. We used the official campaign hashtag #DDLZAN. The volume of data was 123.268 tweets and 572.000 retweets. Given the large number of documents, we filtered the data. From the tweets' database we selected only the most active authors, who published more than 10 tweets in the stream. Moreover, we used the retweets database to select the influencers from the collected authors. We used the interactions between users to select the ones who were more central in the discourse. Briefly, we selected the most cited, quoted and mentioned authors. This interaction-based criterion allowed us to focus on the core of the data stream and study the influencer opinions. The number of influencers selected is 406. We selected these users from our database, and we aggregated their posts (Song et al. 2014). We calculated the interactivity measures for each author: amplification and endorsement. We have operationalized these concepts as quantitative variables. We operatized the amplification as the mean of the post's retweets and the endorsement as the mean of likes of the posts. We created a database with each author as record and published texts, amplification, and endorsements as fields. The posts are the content variable, and the interactivity measures the covariates associated.

## 5. Results

Twitter data are unstructured, so it's necessary to perform some phases of pre-processing for having structured data. There are different steps:

- Tokenized the documents, obtaining a set of distinct strings (tokens)

separated by spaces or punctuation marks;
- Normalized the text, so convert all the letters of the texts into a lower case;
- Removed special characters, punctuations, and numbers from the dataset. Also, special characters, hashtags (i.e. #ddlzan), symbols and stopwords are eliminated;
- Defined a grammatical tagging, which is the process of marking a word in a text as corresponding to a particular part of speech. In this case, we considered the nouns, verbs, and adjectives.

The pre-processing phase returned a database composed by 17353 tokens, 1182 type and 71 documents. The documents terms matrix dimension was 1182x71. Once this is done, STM can be applied to identify latent themes within the textual corpus. To define the optimal number of $k$ topics to extract, we considered the measures of held-out likelihood and semantic coherence (Fig.1). Generally, topic models with excellent fit turn out to be higher on measures of held-out likelihood and semantic coherence (Roberts et al., 2016). So, we selected $k = 5$.
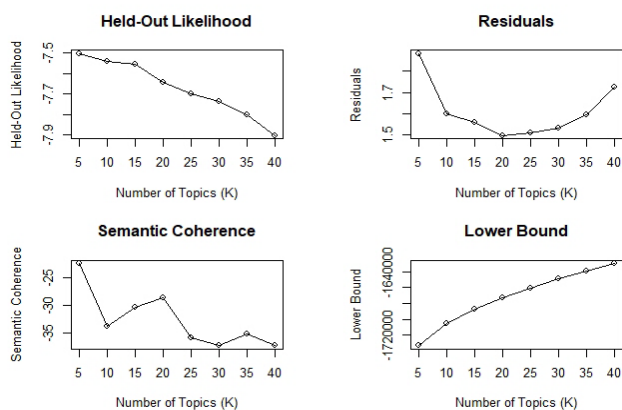


*Figure 1.  Diagnostic values by number of Topics*

In fig.2 there are the top 10 frequent terms associated with the topics. We defined the first topic as "Law principles" because the terms associated with it allow to understand the content of DDLZan, including the fight against "homophobia" and violence motivated by "sexual" orientation and "gender identity".  The second topic concerns the "Propaganda of the right-wing" against the approval of the bill, which has tried to scuttle the law by delaying the vote and asking for continuous changes to DDLZan. In particular, the contrary positions expressed by opposing politicians reiterated concerns already expressed by the Vatican.

Following this, the third topic concerns the criticism made by the center-left parties to Matteo Renzi and his party Italia Viva that, even though they are part of the government majority, have expressed perplexity about the contents of the proposal, asking for an agreement that would allow the law to be approved without including the protection of gender identity. Topic 4 highlights the demands of the majority parties to get to a vote soon, due to continued delays and attempts at ostracism by the opposition. In addition to this, the last topic concerns the democratic and public participation of citizens, who wish to have a law that allows protecting minorities comparing with other European countries.

*Table 1. Top 10 frequent terms associated with topics*

| Topic 1 | Topic 2 | Topic 3 | Topic 4 | Topic 5 |
| --- | --- | --- | --- | --- |
| Libertà | Senato | Votare | Approvare | Parlare |
| Omofobia | Renzi | Senato | Diritto | Italiano |
| Genere | Affossare | Renzi | Chiedere | Pensare |
| Donna | Destra | Chiedere | Parlare | Sinistra |
| Gay | Diritto | Muro | Genere | Approvare |
| Sinistra | Votare | Numero | Votare | Scrivere |
| Identità | Voto | Politico | Identità | Opinione |
| Parlare | Politico | Approvare | Paese | Politico |
| LGBT | Vaticano | Segreto | Odiare | Discriminazione |
| Sesso | Chiedere | Accordo | Anno | Fedez |

## 6. Conclusions

Twitter encourages the formation of both social proximity networks and interest-based ties that can lead to filter bubbles (Marichal et al., 2020). So, we believe that by capturing the most important stakeholders in the debate, we can better understand the dynamics of stakeholder engagement. This happens because more and more individuals, including politicians, actors, bloggers, express their opinions and their daily activities on platforms and influence users by allowing content to spread. Fig. 3 and Fig. 4 plot a Loess, smoothed line of contributions plotted with confidence intervals on the predicted topic proportions for the measure of endorsement and amplification. According to amplification, we visualized that topic 4 has a major value, referring to the need that Twitter users

expressed about having a law against homophobia, while the topic 5 that expressed the participation of citizens has a higher value of endorsement.
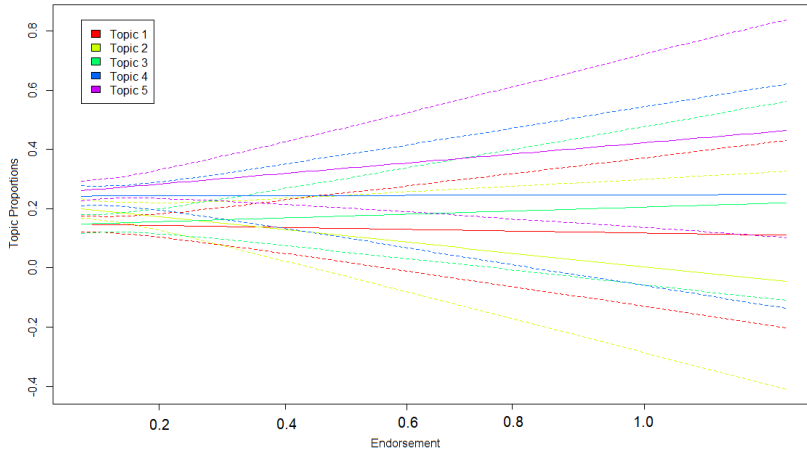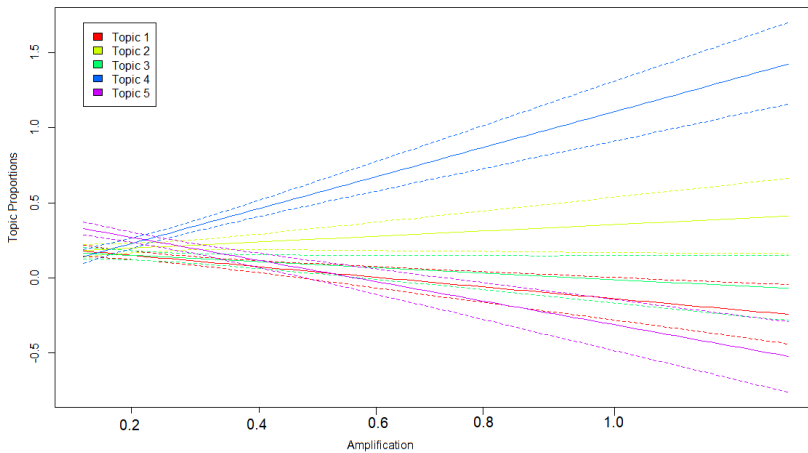


*Figure 2. Topic proportions by endorsement*



*Figure 3. Topic proportions by amplification*

### References

Airoldi E. M., Cohen W. W. and Fienberg S. E. (2004). Bayesian models for frequent terms in text. In *Proceedings of the Classification Society of North America and INTERFACE Annual Meetings*

Blei D.M. and Lafferty J.D. (2007) A correlated topic model of science. *The Annals of Applied Statistics* 1(1):17-35

Castelló I. and Lopez-Berzosa D. (2021). Affects in Online Stakeholder Engagement: A Dissensus Perspective. *Business Ethics Quarterly,* 1-36.

Chatziadam P. Dimitriadis A., Gikas S., Logothetis I., Michalodimitrakis M., Neratzoulakis M. and Kondylakis H. (2020). TwiFly: A Data Analysis Framework for Twitter. *Information, 11(5)*, 247.

Colleoni E. Rozza A. and Arvidsson A. (2014). Echo Chamber or Public Sphere? Predicting Political Orientation and Measuring Political Homophily in Twitter Using Big Data. *Journal of Communication*, 64(2), 317–332.

Gokulakrishnan B., Priyanthan P., Raghavan T., Prasath N., and Perera A., (2012) Opinion mining and sentiment analysis on a Twitter data stream, *International Conference on Advances in ICT for Emerging Regions,* 182-188

Grajzl P. and Murrell P. (2019), Toward understanding 17th century English culture: a structural topic model of Francis Bacon's ideas, J. *Comp. Econ.,* 47 (1)

He L., Han D., Zhou X. and Qu Z. (2020). The voice of drug consumers: online textual review analysis using structural topic model. *International Journal of Environmental Research and Public Health,* 17(10), 3648.

Kwon K. H., Chadha M. and Wang F. (2019). Proximity and networked news public: Structural topic modeling of global Twitter conversations about the 2017 Quebec Mosque shooting. *International Journal of Communication*, 13, 24.

Marichal J. and Neve R. (2019). Antagonistic bias: Developing a typology of agonistic talk on Twitter using gun control networks. *Online Information Review*.

Mishler A., Crabb E. S., Paletz S., Hefright B. and Golonka E. (2015). Using structural topic modeling to detect events and cluster Twitter users in the Ukrainian crisis. *In International conference on human-computer interaction* (pp. 639-644).

Roberts M. E., Stewart B. M. and Airoldi E. M., (2016) A model of text for experimentation in the social sciences. *Journal of the American Statistical Association 111(515)*, pp. 988–1003.

Sobkowicz P., Kaschesky M. and Guillaume B. (2012). Opinion mining in social media: Modeling, simulating, and forecasting political opinions on the web. *Government Information Quarterly*, 29(4), 470–479.

Song G., Ye Y., Du X., Huang X. and Bie S. (2014). Short text classification: A survey. Journal of multimedia, 9(5), 635.

Törnberg, A. and Törnberg, P. (2016). Muslims in social media discourse: Combining topic modeling and critical discourse analysis. *Discourse, Context & Media*, 13, 132-142.