

Approximate Classification with Web Ontologies through Evidential Terminological Trees and Forests

Giuseppe Rizzo*, Nicola Fanizzi**, Claudia d'Amato**, Floriana Esposito

*LACAM – Dipartimento di Informatica — Università degli Studi di Bari "Aldo Moro",
Campus Universitario, Via Orabona 4, 70125 Bari, Italy*

Abstract

In the context of the Semantic Web, assigning individuals to their respective classes is a fundamental reasoning service. It has been shown that, when purely deductive reasoning falls short, this problem can be solved as a prediction task to be accomplished through inductive classification models built upon the statistical evidence elicited from ontological knowledge bases. However also these data-driven alternative classification models may turn out to be inadequate when instances are unevenly distributed over the various targeted classes To cope with this issue, a framework based on logic decision trees and ensemble learning is proposed. The new models integrate the *Dempster-Shafer theory* with learning methods for *terminological decision trees* and *forests*. These enhanced classification models allow to explicitly take into account the underlying uncertainty due to the variety of branches to be followed up to classification leaves (in the context of a single tree) and/or to the different trees within the ensemble model (the forest). In this extended paper, we propose revised versions of the algorithms for learning *Evidential Terminological Decision Trees* and *Random Forests* considering alternative heuristics and additional evidence combination rules with respect to our former preliminary works. A comprehensive and comparative empirical evaluation proves the effectiveness and stability

*Principal corresponding author

**Corresponding author

Email addresses: giuseppe.rizzo1@uniba.it (Giuseppe Rizzo),
nicola.fanizzi@uniba.it (Nicola Fanizzi), claudia.damato@uniba.it (Claudia d'Amato),
floriana.esposito@uniba.it (Floriana Esposito)

of the classification models, especially in the form of ensembles.

Keywords: ontologies, logic decision trees, Dempster-Shafer theory, instance classification

1. Introduction

Sharing knowledge that is encoded along formal ontologies, thus enabling rich reasoning capabilities, plays a key role in the context of the *Semantic Web* (SW). However, standard deductive inference mechanisms sometimes show their limitations because of the inherent incompleteness of the ontological knowledge bases combined with the adoption of an open-world semantics, which is natural in such a Web-scale heterogeneous and distributed context.

In order to tackle the consequences of these distinctive aspects, alternative forms of reasoning, based on statistical models that can be induced through data-driven methods, have been introduced for performing various tasks such as *concept retrieval* and *query answering* [1] more effectively. It has been shown that these tasks have been cast as *classification* problems, which amount to deciding the membership of an individual with respect to a target concept, and they have been solved through inductive learning methods exploiting statistical regularities in the underlying knowledge base. Specifically, the resulting models have been used by approximate classification procedures applied to the knowledge bases also in combination with deductive inference services [2]. The application of these methods has shown interesting results such as the ability to synthesize new concepts and/or produce inductive classification models inspired by *Inductive Logic Programming* (ILP) like *terminological decision trees* [3], i.e. *logic decision trees* [4, 5] whose inner node tests are expressed in terminological languages (that is Description Logics [6]). Additionally, exploiting such statistical models, non logically-derivable yet still consistent assertional knowledge may be suggested.

However, such alternative methods and models have also revealed some shortcomings. One of the issues is that they do not allow an explicit repre-

27 sentation of uncertainty to be specifically exploited for managing those cases
28 when the classification procedure assigns an uncertain membership. To better
29 tackle these cases, an enhanced model, called *evidential terminological decision*
30 *tree* has been devised, by integrating primitives of the *Dempster-Shafer The-*
31 *ory* [7]. The main advance with respect to terminological decision trees regards
32 the heuristic used to select the concept installed into inner nodes (based on the
33 *non-specificity measure* [8] rather than the classic measures stemming from *in-*
34 *formation gain*) and the classification procedure (that explores all the possible
35 paths departing from a node with an uncertain test result).

36 Another issue concerns the distribution of the training data. In general, the
37 individuals that are known (or can be logically assessed as) positive and negative
38 instances for a given target concept (that is those that are instances of a target
39 concept or of the negated target concept) may not be equally distributed. This
40 skewness may be noticeably larger when considering individuals whose mem-
41 bership cannot be assessed by reasoning under an open-world semantics. This
42 class-imbalanced setting may affect the model, resulting in poor performances.
43 Various methods have been devised to tackle the general unbalance learning
44 problem (see [9] for a survey of the various approaches). As regards the specific
45 task of learning instance classification models for inductive query answering on
46 SW knowledge bases, we investigated the adoption of methods for *ensemble*
47 *models* [10] that are made up of a certain number of classifiers, trained by the
48 so-called *weak learners*, and whose final prediction results from the combination
49 of the predictions made by each classifier. Specifically, the combination is given
50 by a specific rule playing the role of the *meta-learner*. Particularly, we proposed
51 an algorithm for inducing *terminological random forests* [10] that extends (*First*
52 *Order*) *random forests* [11, 12] with the use of Description Logics: the model is
53 an ensemble of terminological decision trees [3].

54 Employing these models, the membership of a test individual w.r.t. a target
55 concept is decided according to a majority vote rule (although various other
56 strategies for combining predictions have been proposed [13, 14, 15]): each clas-
57 sifier equally contributes to the final decision returning a vote in favor of a

58 single membership. In this way, some other aspects are not considered explic-
59 itly, such as the uncertainty about the single membership-label assignments
60 and the disagreement that may intervene among weak learners. Particularly,
61 the latter issue is crucial for the performance of ensemble models [16]: using the
62 aforementioned type of forests, we noted that most misclassification cases were
63 related to situations in which votes are evenly distributed with respect to the
64 admissible labels. A weighted voting procedure may be an alternative strategy
65 to mitigate the problem, but it requires a criterion for setting the weights.

66 In this sense, introducing a meta-learner which can manipulate the *soft* pre-
67 dictions made by each classifier (i.e. a prediction with a confidence measure for
68 each membership value) rather than *hard predictions* (where only the predicted
69 label is returned) may be a solution. Adopting the random forests as ensem-
70 bles, this can be accomplished by considering evidential terminological decision
71 trees [7] as base models. Dempster-Shafer theory has already been used in
72 combination with ensemble learning procedures (e.g. see [17]). However, most
73 of the methods apply to problems that involve simpler knowledge representa-
74 tions. Additionally, none of them has been employed for predicting assertions
75 on ontological knowledge bases.

76 Therefore, we further extended the model proposing a framework for the in-
77 duction of *Evidential Terminological Random Forests* for ontological knowledge
78 bases [18]. Employing evidential terminological decision trees, the approach
79 does not require the computation of decision templates. After the induction
80 of the forest, new individuals are classified by combining the evidence on the
81 membership prediction made by each tree through Dempster’s rule [19].

82 However, we noted that the proposed framework had some limitations [7, 18].
83 Firstly, the heuristic to select the most promising label adopted by evidential
84 terminological decision tree learning algorithm did not consider the presence of
85 conflicting evidence. Secondly, the combination rule represented a bottleneck
86 of the classification step: therefore it is important to investigate alternative
87 solutions for improving the efficiency of the classification. Thirdly, the size of
88 evidential terminological random forests seemed not to affect the predictive-

89 ness of the ensemble model (due to a weak diversification of the ensemble) but
90 represented a source of complexity during the classification step.

91 Consequently, in this paper we extended the framework for learning ev-
92 idential terminological decision trees and random forests along the following
93 directions:

- 94 • we used different heuristics based on other total uncertainty measures
95 (than the sole non-specificity measure) to drive the selection of the con-
96 cepts to be installed into the nodes of evidential terminological decision
97 trees;
- 98 • we used further combination rules to pool the evidence obtained by
99 traversing each tree;
- 100 • we used further combination rules as meta-learner for evidential termino-
101 logical random forests;
- 102 • we set up a comprehensive and comparative experimental evaluation show-
103 ing the effectiveness of the proposed extensions when performing inductive
104 instance retrieval.

105 The remainder of the paper is organized as follows: the next section intro-
106 duces basics on the targeted representation language and the problem we aim
107 to solve, that is inducing classifiers for the SW context; Sect. 3 recalls the ba-
108 sics on Dempster-Shafer Theory, required for understanding the framework of
109 the evidential tree-based models presented in Sect. 4. In Sect. 5, the empirical
110 evaluation of the classification models is described, while Sect. 6 discusses re-
111 lated approaches. Sect. 7 draws conclusions and illustrates some perspectives
112 for further developments.

113 2. Basics of Description Logics and Problem Definition

114 In this section we recall the basics of *Description Logics* (DLs), that is the
115 family of knowledge representation languages at the core of the standard *Web*

116 *ontology language*¹ (OWL - DL).

117 In DLs, a domain is modeled in terms of a set of *atomic concepts*, $N_C =$
 118 $\{A, B, \dots\}$ and *atomic roles*, $N_R = \{R, S, \dots\}$. Two noteworthy concepts are
 119 the *top concept*, denoted with \top , and the *bottom concept*, denoted with \perp . DLs
 120 are endowed with a set of operators to combine atomic concepts and forming
 121 complex descriptions, such as complement, conjunction and disjunction. A set
 122 of constants, dubbed as *individuals* and denoted with $N_I = \{a, b, \dots\}$, is to be
 123 considered as the names of the objects of the domain to be represented.

124 The semantics of the constructs is defined in terms of *interpretations*. An
 125 interpretation is a couple $\mathcal{I} = (\Delta, \cdot^{\mathcal{I}})$, where its *domain* $\Delta^{\mathcal{I}}$ is a non-empty set of
 126 objects while $\cdot^{\mathcal{I}}$ is the *interpretation function* that maps each concept $C \in N_C$
 127 onto a set of objects $C^{\mathcal{I}} \subseteq \Delta^{\mathcal{I}}$ and each role $R \in N_R$ onto a binary relation
 128 $R^{\mathcal{I}} \subseteq \Delta^{\mathcal{I}} \times \Delta^{\mathcal{I}}$. In addition, $\top^{\mathcal{I}} = \Delta^{\mathcal{I}}$ and $\perp^{\mathcal{I}} = \emptyset$. The semantics of complex
 129 concept descriptions is defined recursively depending on the available operators
 130 for building complex concepts. For instance, for the case of \mathcal{ALC} , the semantics
 131 of complex description is defined as follows:

- 132 • $(D \sqcap E)^{\mathcal{I}} = D^{\mathcal{I}} \cap E^{\mathcal{I}}$
- 133 • $(\neg D)^{\mathcal{I}} = \Delta^{\mathcal{I}} \setminus D^{\mathcal{I}}$
- 134 • $(\forall R.D)^{\mathcal{I}} = \{a \in \Delta^{\mathcal{I}} \mid \forall b \in \Delta^{\mathcal{I}}, (a, b) \in R^{\mathcal{I}} \rightarrow b \in D^{\mathcal{I}}\}$
- 135 • $(\exists R.D)^{\mathcal{I}} = \{a \in \Delta^{\mathcal{I}} \mid \exists b \in \Delta^{\mathcal{I}}, (a, b) \in R^{\mathcal{I}} \wedge b \in D^{\mathcal{I}}\}$

136 Finally, each individual name is mapped onto an element of $\Delta^{\mathcal{I}}$.

137 A knowledge base is a pair $\mathcal{K} = (\mathcal{T}, \mathcal{A})$ where \mathcal{T} and \mathcal{A} denote its TBox and
 138 ABox. The TBox contains intensional knowledge about the domain, modeled
 139 as inclusion axioms $C \sqsubseteq D$ (meaning that D subsumes C) and interpreted
 140 as $C^{\mathcal{I}} \subseteq D^{\mathcal{I}}$ for every interpretations \mathcal{I} . Given two concepts C and D , C is
 141 equivalent to D if for every interpretations \mathcal{I} , $C^{\mathcal{I}} = D^{\mathcal{I}}$. Alternatively, C and D
 142 are equivalent if $C \sqsubseteq D$ and $D \sqsubseteq C$. The ABox \mathcal{A} contains factual knowledge,

¹www.w3.org/OWL

143 i.e. assertions concerning individuals. In the ABox there are two kinds of
 144 assertions: *concept* $C(a)$ and *role assertions* $R(a, b)$. The set of individuals
 145 occurring in \mathcal{A} are denoted by $\text{Ind}(\mathcal{A})$.

146 Knowledge bases are also equipped with deductive reasoning capabilities.
 147 An important reasoning service for our purposes is *instance checking*: an indi-
 148 vidual a is an instance of a concept C if, for every model of \mathcal{K} , $C(a)$ holds. This
 149 can be denoted with $\mathcal{K} \models C(a)$. We will be also interested in the case where
 150 $\mathcal{K} \models \neg C(a)$. These instances will be exploited as *examples* (positive and neg-
 151 ative examples respectively) in our learning procedures. Note that, due to the
 152 reasoning under the *Open World Assumption* (OWA) that is generally adopted
 153 in this context, it may happen that $C(a)$ and $\neg C(a)$ are satisfied by different
 154 models of \mathcal{K} . This means that neither $\mathcal{K} \models C(a)$ nor $\mathcal{K} \models \neg C(a)$ holds, i.e.
 155 there is insufficient knowledge to decide the membership of a w.r.t. the target
 156 concept using standard deductive inference services. Such individuals will be
 157 considered as instances with *uncertain membership* w.r.t. C .

158 In order to overcome this inherent limitation, it is possible to resort to de-
 159 cision procedures that are based on inductive (statistical) classification models.
 160 They can be learned by fitting a function from available examples (individuals
 161 for which the membership w.r.t. C is known) that amounts to solving a mini-
 162 mization problem based on a notion of misclassification *risk*. A general learning
 163 task aiming at classification models can be defined as follows:

164 **Definition 1 (learning problem).**

165 *Given*

- 166 • a target concept C
- 167 • a set of instances \mathbf{E}
- 168 • a set of labels \mathcal{L} to denote the membership w.r.t. C
- 169 • a joint probability distribution between \mathbf{E} and \mathcal{L} , namely $P(\mathbf{E}, \mathcal{L})$, mea-
 170 suring the chance of an element of \mathbf{E} to be assigned with one of the labels

- 171 • a set of hypotheses $\mathcal{H} = \{h : \mathbf{E} \rightarrow \mathcal{L}\}$, i.e. classification functions that
 172 can predict a label for their arguments
- 173 • a loss function $L : \mathbf{E} \times \mathcal{L} \rightarrow [0, +\infty[$ to assign a penalty for predicting an
 174 incorrect label for a given instance

175 **Find** a function $h^* \in \mathcal{H}$ such that:

$$h^* = \arg \min_{h \in \mathcal{H}} \mathbb{E}_P [L(h(a), l)] \quad (1)$$

176 This definition requires the expected risk to be computed over the data gen-
 177 erating distribution P , which is usually unknown. Therefore a more concrete
 178 definition will be specified for the case of DL knowledge bases, aimed at inducing
 179 a classification function that minimizes an *empirical risk* of error on the training
 180 set, and it can be reformulated for the targeted representation as follows:

181 **Definition 2 (learning classifiers for DL knowledge bases).**

182 *Given*

- 183 • a target concept C in the signature of a knowledge base $\mathcal{K} = \langle \mathcal{T}, \mathcal{A} \rangle$
- 184 • a set of membership labels $\mathcal{L} = \{-1, 0, +1\}$ to denote, resp., the positive,
 185 uncertain and negative membership w.r.t. C
- 186 • a loss function $L : \text{Ind}(\mathcal{A}) \times \mathcal{L} \rightarrow [0, +\infty[$
- 187 • a training set of examples for which the correct labels are known, i.e. the
 188 values of a correct classifier² $f : \text{Ind}(\mathcal{A}) \rightarrow \mathcal{L}$, $\mathbf{Tr} = \mathbf{P} \cup \mathbf{N} \cup \mathbf{U}$ where:
- 189 $\mathbf{P} = \{a \in \text{Ind}(\mathcal{A}) \mid f(a) = +1\}$ i.e. $\{a \in \text{Ind}(\mathcal{A}) \mid \mathcal{K} \models C(a)\}$
- 190 $\mathbf{N} = \{a \in \text{Ind}(\mathcal{A}) \mid f(a) = -1\}$ i.e. $\{a \in \text{Ind}(\mathcal{A}) \mid \mathcal{K} \models \neg C(a)\}$
- 191 $\mathbf{U} = \{a \in \text{Ind}(\mathcal{A}) \mid f(a) = 0\}$ i.e. $\{a \in \text{Ind}(\mathcal{A}) \mid \mathcal{K} \not\models C(a) \wedge \mathcal{K} \not\models \neg C(a)\}$
- 192 • a set of classification functions, or hypotheses, $\mathcal{H} = \{h : \text{Ind}(\mathcal{A}) \rightarrow \mathcal{L}\}$

²i.e. whose analytic form is not available.

193 **Find** a classification function $h^* \in \mathcal{H}$, approximating f , such that

$$h^* = \arg \min_{h \in \mathcal{H}} \frac{1}{|\mathbf{Tr}|} \sum_{a \in \mathbf{Tr}} L[h^*(a), f(a)] \quad (2)$$

194 Note that the hypothesis set \mathcal{H} acts as a form of *bias* and can be properly defined
 195 in order to exclude trivial solutions (overfitting) such as classifiers induced by
 196 a *rote learner* based on functions that merely memorize the correct labeling for
 197 the training examples (that would be equivalently described by the disjunction
 198 of very specific concepts – one per positive example). Conversely, the aim is to
 199 obtain a solution that is able to ensure a good generalization, that is the ability
 200 to correctly predict the membership for *unseen* individuals, i.e. individuals that
 201 have not been considered during the training phase. In this paper, we present
 202 a solution to this learning problem based on a tree classification model which
 203 combines logics and evidence-based prediction.

204 3. Basics of the Dempster-Shafer Theory

205 In this section basics of Dempster-Shafer Theory are summarized since it
 206 represents the main building block for the formalization of the evidential tree-
 207 based models presented in Sect. 4.

208 The Dempster-Shafer Theory (DST) [20] can be regarded as a generalization
 209 of the *Bayesian subjective probability* theory. The framework offers an alterna-
 210 tive to traditional probabilistic theory for the mathematical representation of
 211 uncertainty: a probability mass can be assigned to a set or an interval without
 212 knowing the probability of the specific elements. As argued in [20], this aspect
 213 may be a valuable tool when knowledge is obtained from expert elicitation.

214 In the DST, the *frame of discernment* is a set of exhaustive and mutually
 215 exclusive hypotheses $\Omega = \{\omega_1, \omega_2, \dots, \omega_n\}$ about a domain. Moving from the
 216 frame of discernment we can defined the *basic belief assignment*

217 **Definition 3 (BBA and focal element).** A basic belief assignment (*BBA*)
 218 is defined as a mapping where $m : 2^\Omega \rightarrow [0, 1]$ so that $m(\emptyset) = 0$ and
 219 $\sum_{A \in 2^\Omega} m(A) = 1$. If $m(A) > 0$, A is a focal element for m .

220 Note that BBAs extend the standard probability measures. The key dif-
 221 ference regards the relaxation of the *monotonicity* property of the probability
 222 measures. Indeed the degrees of belief are ascribed to sets of events rather than
 223 to single events. This means that for a BBA, given $A, B \in 2^\Omega$, $A \subseteq B$ does not
 224 imply $m(A) \leq m(B)$. This property derives from committing the value $m(B)$
 225 only to the set B and not to any of its subsets. Conversely, in the probability
 226 theory, the probability of an event A is exactly the sum of probabilities assigned
 227 to the single $a \in A$.

228 Other functions can be derived from BBAs such as *belief* and *plausibility*.

Definition 4 (belief). *The belief in A , denoted by $\text{Bel}(A)$, represents a mea-
 sure of the support committed to A given the available evidence:*

$$\text{Bel}(A) = \sum_{B \subseteq A} m(B).$$

Definition 5 (plausibility). *The plausibility of A , $\text{Pl}(A)$, represents the total
 belief that may be committed to A when further evidence becomes available:*

$$\text{Pl}(A) = \sum_{B \cap A \neq \emptyset} m(B).$$

229 Note that, differently from the BBA m , Bel and Pl are monotonic. As described
 230 in the following, this is taken into account when these measures are used with
 231 the models proposed in this paper.

232 3.1. Combination Rules

233 Combination rules are operators for pooling information obtained from mul-
 234 tiple sources. These sources provide different assessments for the frame of
 235 discernment of the domain of interest. DST traditionally assumes that these
 236 sources are independent, although this constraint has been progressively relaxed
 237 with the introduction of new rules.

238 Many operations have been proposed in the literature [19]. In the sequel,
 239 we briefly survey the most important combination rules. In the rest of the
 240 paper, we will denote the application of one of such combination rules on two
 241 (or multiple) BBAs with the symbol \oplus (e.g. $m_{1,2} = m_1 \oplus m_2$).

242 *3.1.1. Dempster's Rule*

243 The original combination rule of multiple BBAs known as Dempster's rule is
 244 a generalization of Bayes' rule [21]. The resulting BBA can be computed with:

$$\forall A \subseteq \Omega \quad m_{1,2}(A) = \begin{cases} \frac{1}{1-c} \sum_{B \cap C = A} m_1(B) \cdot m_2(C) & \text{if } A \neq \emptyset \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

where

$$c = \sum_{B \cap C = \emptyset} m_1(B) m_2(C) \quad (4)$$

245 This rule emphasizes the agreement between the sources adopting the normal-
 246 izing factor c to distribute the conflicting evidence. It has come under serious
 247 criticism when the amount of conflict among sources is significant leading to
 248 counterintuitive results.

Example 1 (Dempster's rule). *Let us consider two BBAs m_1 and m_2 defined over a simple frame of discernment $\Omega = \{\omega_1, \omega_2, \omega_3\}$ whose focal elements are reported below:*

$$m_1(\{\omega_1\}) = 0.99, \quad m_1(\{\omega_3\}) = 0.01, \quad m_2(\{\omega_2\}) = 0.99, \quad m_2(\{\omega_3\}) = 0.01$$

Applying the rule, the pooled BBA value for $\{\omega_3\}$ is:

$$m_{1,2}(\{\omega_3\}) = \frac{1}{1 - 0.99 \cdot 0.99 - 0.99 \cdot 0.01 - 0.99 \cdot 0.01} \cdot 0.01 \cdot 0.01 = 1$$

249 *Note that this result is due to the agreement of the evidence in favor of $\{\omega_3\}$*
 250 *and the disagreement between $\{\omega_1\}$ and $\{\omega_2\}$.*

251 To prevent cases like the one reported above, which may affect the effectiveness
 252 of the models described in this paper, we investigated the effectiveness of further
 253 rules.

254 *3.1.2. Dubois-Prade Disjunctive Pooling Rule*

255 This rule [22] takes into account the union of the probability masses (dis-
 256 junctive rule): this prevents the generation of conflict as there is no rejection of

257 information coming from the various sources. The combination rule is defined
 258 as follows:

$$\forall A \subseteq \Omega \quad m_{1,2}(A) = \sum_{B \cup C = A} m_1(B)m_2(C). \quad (5)$$

259 The union does not generate any conflict and does not reject any informa-
 260 tion asserted by the sources. As such, no normalization is required (unlike the
 261 Dempster’s rule). The drawback of this rule is that it may yield a more impre-
 262 cise result than desirable. It is easy to see that this rule is commutative and
 263 associative.

264 3.1.3. *Mixing*

265 This rule (also known as *averaging*) represents an extension of the average
 266 for probability distributions computed on the BBAs and describes the frequency
 267 of the various values within a range of possible values. The resulting BBA can
 268 be obtained merely as a weighted average of the masses according to the various
 269 features:

$$\forall A \subseteq \Omega \quad m_{1,\dots,n}(A) = \frac{1}{n} w_i m(A) \quad (6)$$

270 where a normalized weight vector \mathbf{w} is generally considered. The values of the
 271 weights should reflect a degree of confidence in the sources. This rule is commu-
 272 tative, idempotent and quasi-associative³. For our purpose, we are interested
 273 in associative combination rules to prevent the final decision to be affected by
 274 the pooling order of the considered BBAs, namely *Dempster’s rule* and *Dubois-*
 275 *Prade rule*. In the experiments we will consider also mixing rule to investigate

³ A quasi-associative operation is an operation that can be broken down in two associative sub-operations. For instance, the mathematical average is quasi-associative: the value is obtained as the sum of a list of numbers divided by the number of the elements in the list (both the sum of the terms and the counting of the element in the list are associative operations)

276 the effectiveness of the predictive models when a quasi-associative rule is em-
 277 ployed.

278 3.2. Measures of Total Uncertainty

279 In the context of the DST, various measures of uncertainty can be considered.
 280 These measures are typically defined as generalizations of *Shannon's entropy*
 281 or of other types of measures of uncertainty proposed in Probability Theory.
 282 Alternatively, they can be determined according to the conflict existing among
 283 the BBAs to be pooled according to a given combination rule. In this section,
 284 we briefly recall some measures. For more details, see [8].

285 The *non-specificity* measure [23] quantifies the degree of imprecision related
 286 to a BBA:

$$NS(m) = \sum_{A \in 2^\Omega} m(A) \log(|A|) \quad (7)$$

287 The measure of *confusion* is defined on the ground of a BBA and the belief
 288 measure, as reported below [24]:

$$Confusion(m) = - \sum_{A \in 2^\Omega} m(A) \log(Bel(A)) \quad (8)$$

289 The measure of *dissonance* [25] is based on a BBA and the plausibility and
 290 is defined as follows:

$$Dissonance(m) = - \sum_{A \in 2^\Omega} m(A) \log(Pl(A)) \quad (9)$$

291 In the sequel, we will adopt these criteria to select the best features that
 292 compose the model proposed in this paper.

293 4. Evidence-based Terminological Trees and Forests

294 The original notions of terminological decision trees and random forests will
 295 be now recalled before introducing the new methods for the induction and usage
 296 of the evidence-based versions of these classification models.

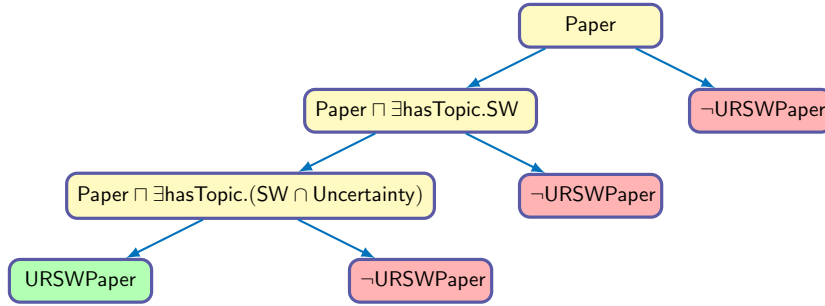


Figure 1: A TDT for predicting if a paper may have appeared in URSW proceedings

297 4.1. Terminological Decision Trees and Random Forests

298 Classification can be performed by inducing *terminological decision trees*
 299 (TDTs) [3]. A TDT is basically a binary tree whose leaves contain labels that
 300 denote the (positive/negative) membership with respect to the target concept;
 301 each inner node, dubbed also *decision* or *test* node, contains a DL concept
 302 description D (in conjunctive form) and the descending edges from such a node
 303 represent the result of a test over D (positive, negative).

304 Fig. 1 illustrates a simple example of a TDT that can describe the individ-
 305 uals of a knowledge base that are papers appeared in the URSW proceedings
 306 (target concept `URSWPaper`). Note that, given a node with a concept descrip-
 307 tion D , its left child may be either a leaf or another decision node containing a
 308 concept description E such that $E \sqsubseteq D$, whereas the right child may be either
 309 another leaf or a decision node containing a concept description E' for which
 310 $E' \sqsubseteq \neg D$ is intended. For instance, the root node contains the concept `Paper`
 311 while its left child is another decision node containing the concept description
 312 `Paper \sqcap \exists\text{hasTopic.SW}` and its right child is a leaf with a negative label.

313 Similarly to other supervised models, the predictiveness of a TDT can be
 314 affected by the *class-imbalance problem*. In Machine Learning, this problem con-
 315 cerns the skewness of the training data distributions. Especially in multi-label
 316 settings, the problem occurs when the number of training examples belonging
 317 to a particular category (the *majority class*) overwhelms the number of those
 318 belonging to the others.

319 In order to tackle this problem, the most common approaches that have been
320 proposed are based on a sampling strategy [26]. One of the simplest methods is
321 an *under-sampling* strategy that randomly discards instances belonging to the
322 majority class in order to re-balance the dataset. However, this method causes
323 a loss of information due to the possible removal of useful (*critical*) examples
324 that may be essential for inducing a predictive model.

325 A *terminological random forest* (TRF) is an ensemble model trained through
326 a procedure that combines a random under-sampling strategy with ensemble
327 learning [10]. A TRF is basically made up of a certain number of TDTs, where
328 each of them is built by considering a (quasi-)balanced dataset. The ensemble
329 model assigns the final classification for a new individual by appealing to a
330 majority vote procedure. Therefore each TDT returns a *crisp prediction*: each
331 provides an equal contribution to the final decision regarding the membership
332 label, as no measure of confidence is available per single prediction.

333 In order to consider also this kind of information and tackling other relevant
334 problems related to the uncertainty about the class assignments (e.g. cases of
335 ties in conflicting predictions) and the disagreement between classifiers that may
336 lead to misclassifications [10], we need to resort to other models for the ensemble
337 approach.

338 4.2. Evidential Terminological Decision Trees

339 To better take into account the mentioned forms of uncertainty, it has been
340 shown how approximate class-membership prediction can be carried out by in-
341 ducing *evidential terminological decision trees* (ETDTs) [7], an extension of the
342 TDTs based on the DST. ETDTs are defined in a similar way with respect to
343 TDTs. However, unlike TDTs, each inner node contains a pair $\langle D, m \rangle$ where,
344 besides the concept description D , there is a BBA m based on the membership
345 of the individuals w.r.t. D .

346 Fig. 2 reports an example of ETDT used for deciding whether a paper has
347 been published in the proceedings of URSW. Similarly to a TDT, each decision
348 node contains a concept description D , while the left (resp. right) child may

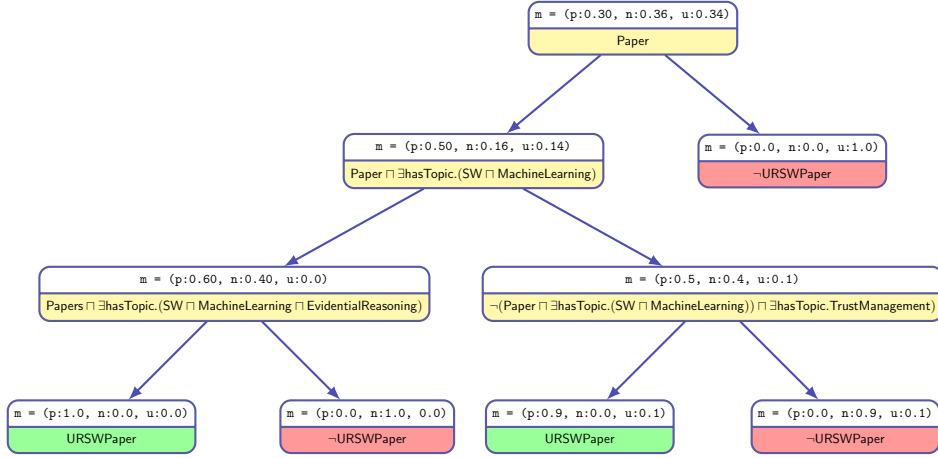


Figure 2: An ETDT for deciding if a paper has appeared in the URSW proceedings

349 either be a leaf (containing the corresponding label) or another decision node
 350 with a concept description $E \sqsubseteq D$ (resp. $E \sqsubseteq \neg D$). In addition, each node
 351 contains also a BBA, which can be estimated from the training instances used
 352 to learn the model, as described in the sequel.

353 4.2.1. Growing ETDTs

354 Before presenting the learning procedure we need to introduce the some
 355 notation. Moving from the formulation of the learning problem [10] (defined
 356 in Sect. 2), we will use the subset of the definite classification labels, $\Omega =$
 357 $\{-1, +1\} \subset \mathcal{L}$, as the *frame of discernment* of the problem (see Sect. 3).
 358 Therefore the positive membership label $+1$ corresponds to the subset $\{+1\}$
 359 of the frame of discernment, the negative membership label -1 corresponds to
 360 the subset $\{-1\}$, and the case of uncertain-membership will be denoted with the
 361 label 0 corresponding to $\{-1, +1\}$.

362 Practically, to learn an ETDT model, a *divide-and-conquer* approach is
 363 adopted where a set of (more specific) concept descriptions is generated from
 364 the one contained in parent nodes. For each specialization, a BBA is also com-
 365 puted. Then the best description (and the corresponding BBA) is selected, e.g.

Algorithm 1 The routines for inducing ETDTs

```

1 const  $\theta \in ]0, 1]$  {min.\ purity threshold parameter}
2
3 function INDUCEETDTREE( $\langle \mathbf{P}, \mathbf{N}, \mathbf{U} \rangle, C, D, m, \hat{\text{Pr}}$ ):  $T$ 
4 input  $\langle \mathbf{P}, \mathbf{N}, \mathbf{U} \rangle$ : training set;  $C$ : target concept;  $D$ : concept,  $m$ : BBA;  $\hat{\text{Pr}}$ : priors
5 output  $T$ : ETDT
6 begin
7  $T \leftarrow$  new ETDT
8 if  $|\mathbf{P}| = 0$  and  $|\mathbf{N}| = 0$  then
9   if  $\hat{\text{Pr}}(+1) \geq \hat{\text{Pr}}(-1)$  then {pre-defined constants wrt the whole training set}
10     $T.\text{root} \leftarrow \langle C, m \rangle$ 
11   else
12     $T.\text{root} \leftarrow \langle -C, m \rangle$ 
13 else if  $(m(\{-1\}) \simeq 0)$  and  $(m(\{+1\}) > \theta)$  then
14    $T.\text{root} \leftarrow \langle C, m \rangle$ 
15 else if  $(m(\{+1\}) \simeq 0)$  and  $(m(\{-1\}) > \theta)$  then
16    $T.\text{root} \leftarrow \langle -C, m \rangle$ 
17 else
18    $\mathbf{S} \leftarrow \emptyset$ 
19   for  $E \in \rho(D)$  {assignBBA for each candidate}
20     $m_E \leftarrow$  COMPUTEBBA( $E, \langle \mathbf{P}, \mathbf{N}, \mathbf{U} \rangle$ )
21     $\mathbf{S} \leftarrow \mathbf{S} \cup \{ \langle E, m_E \rangle \}$ 
22    $\langle E^*, m^* \rangle \leftarrow$  SELECTBESTCANDIDATE( $\mathbf{S}$ )
23    $\langle \langle \mathbf{P}^l, \mathbf{N}^l, \mathbf{U}^l \rangle, \langle \mathbf{P}^r, \mathbf{N}^r, \mathbf{U}^r \rangle \rangle \leftarrow$  SPLIT( $E^*, \langle \mathbf{P}, \mathbf{N}, \mathbf{U} \rangle$ )
24    $T.\text{root} \leftarrow \langle E^*, m^* \rangle$ 
25    $T.\text{left} \leftarrow$  INDUCEETDT( $\langle \mathbf{P}^l, \mathbf{N}^l, \mathbf{U}^l \rangle, C, E^*, m^*, \hat{\text{Pr}}$ )
26    $T.\text{right} \leftarrow$  INDUCEETDT( $\langle \mathbf{P}^r, \mathbf{N}^r, \mathbf{U}^r \rangle, C, \neg E^*, m^*, \hat{\text{Pr}}$ )
27 return  $T$ 
28 end

```

366 the one having the smallest non-specificity w.r.t. the previous level.

367 Alg. 1 illustrates the training procedure. It distinguishes various cases: the
 368 non-recursive ones are those for which leaves are defined while the final one
 369 determines the inner nodes, hence the subtree structure, recursively.

370 The first case copes with the lack of examples ($|\mathbf{P}| = 0$ and $|\mathbf{N}| = 0$) routed
 371 to the node resorting to the prior probability (estimates).

372 The following cases determine the label for a leaf-node when it is (sufficiently)
 373 pure, i.e. no positive (resp. negative) example is found (or just a few) while most
 374 of the examples are negative (resp. positive). This *purity* condition is evaluated
 375 by considering the BBA m given as an input to the algorithm ($m(\{-1\}) \simeq 0$ and
 376 $m(\{+1\}) > \theta$ or $m(\{+1\}) \simeq 0$ and $m(\{-1\}) > \theta$), where θ is a purity threshold.
 377 The values of a BBA function for the membership values are obtained from the
 378 distribution of positive, negative and uncertain-membership instances w.r.t. the
 379 current concept.

380 Finally, the last (recursive) case concerns the availability of a nonnegligible

381 number of both negative and positive examples. In this case, the current concept
 382 description D has to be specialized by means of an operator exploring the search
 383 space of downward refinements of D . Following the approach described in [10,
 384 12], the refinement step produces a set of candidate specializations $\rho(D)$. A
 385 BBA m_E is then built for each candidate $E \in \rho(D)$. Again, the function
 386 can be obtained by counting the number of positive, negative and uncertain-
 387 membership instances). Then the best pair $\langle E^*, m^* \rangle \in \mathbf{S}$ according to the non-
 388 specificity measure is determined by the SELECTBESTCANDIDATE procedure and
 389 finally installed in the current node. Specifically, the procedure tries to find the
 390 pair $\langle E^*, m^* \rangle$ having the smallest non-specificity measure. As an alternative, the
 391 best concept description can be selected in order to maximize either confusion
 392 measure or the dissonance measure w.r.t. the previous level.

393 After the assessment of the best test concept description E^* , the individuals
 394 are partitioned by the procedure SPLIT for the left or right branch according
 395 to the result of the test w.r.t. E^* , maintaining the same group⁴ ($\mathbf{P}^{l/r}, \mathbf{N}^{l/r}$,
 396 or $\mathbf{U}^{l/r}$). Note that a training example a is replicated in both children in
 397 case both $\mathcal{K} \not\models E^*(a)$ and $\mathcal{K} \not\models \neg E^*(a)$ (test with a non-definite, positive or
 398 negative, outcome). The divide-and-conquer strategy is applied recursively until
 399 the instances routed to a node satisfy one of the stopping conditions discussed
 400 above.

401 From a *learning-as-search* perspective, one may regard the induction of an
 402 ETDT as a search process in a hypothesis space \mathcal{H} defined by the set of all
 403 possible ETDTs ruling out those having a sole inner node in the form of a pair
 404 (\top, m) .

405 4.2.2. Prediction

406 Given a test individual a and the induced ETDT, the membership can be
 407 assessed by following one or more paths in the tree. The procedure is reported

⁴Note that the group is related to the membership w.r.t. the target class, while the branch direction depends on the outcome of the test w.r.t. E^* .

Algorithm 2 Class-membership prediction routine through ETDT

```
1 const  $\varepsilon \in ]0, 1]$  {decision threshold parameter}
2
3 function CLASSIFYBYETDT( $a, T$ ) :  $l$ 
4 input  $a$ : individual;  $T$ : ETDT
5 output  $l \in \mathcal{L}$ 
6 begin
7  $M \leftarrow \text{GETLEAFBBALIST}(a, T)$  {list of BBAs located at leaf-nodes}
8  $\bar{m} \leftarrow \bigoplus_{m \in M} m$ 
9 for each  $\emptyset \neq s \in 2^\Omega$  do
10   Compute  $\overline{Bel}(s)$  from  $\bar{m}$ 
11 if  $|\overline{Bel}(\{-1\}) - \overline{Bel}(\{+1\})| \leq \varepsilon$  then
12   predlabel  $\leftarrow 0$  {case of uncertain membership}
13 else
14   predlabel  $\leftarrow \arg \max_{l \in \Omega} \overline{Bel}(\{l\})$  {cases of definite membership}
15 return predlabel
16 end
```

408 in Alg. 2.

409 Specifically, the algorithm traverses recursively the ETDT by performing a
410 test w.r.t. the concept contained in each node that is reached: let $a \in \text{Ind}(\mathcal{A})$
411 and D the concept installed in the current node, if $\mathcal{K} \models D(a)$ (resp. $\mathcal{K} \models \neg D(a)$)
412 the left (resp. right) branch is followed. If neither $\mathcal{K} \models D(a)$ nor $\mathcal{K} \models \neg D(a)$ is
413 verified, both branches are followed.

414 After the exploration of an ETDT (via GETLEAFBBALIST), the list M
415 likely contains multiple BBAs. In this case, the BBAs are pooled according to
416 a combination rule (see Sect. 3) producing \bar{m} .

417 The final decision about the membership to be assigned to the test indi-
418 vidual is made by computing the belief measures for the positive, negative and
419 uncertain membership cases based on the pooled BBA. If the measures for the
420 definite cases are approximately equal (their difference is below a given thresh-
421 old ε), the algorithm will assign the uncertain membership label 0. Conversely,
422 the algorithm selects the definite label ($l \in \Omega$) with higher belief.

423 4.3. Evidential Terminological Random Forests

424 An *evidential terminological random forest* (ETRF) is an ensemble of ET-
425 DTs. We will focus on the procedures for producing an ETRF and for predicting
426 class-membership of input individuals exploiting an ETRF.

Algorithm 3 The routines for inducing ETRFs

```
1 const  $\theta \in ]0, 1]$  {min.\ purity threshold parameter}
2
3 function INDUCEETRF( $\mathbf{Tr}, C, n$ ):  $\mathbf{F}$ 
4 input  $\mathbf{Tr}$  : training set;  $C$  : target concept;  $n \in \mathbb{N}$ 
5 output  $\mathbf{F}$ : ETRF
6 begin
7  $\hat{\mathbf{P}}_{\mathbf{r}} \leftarrow \text{ESTIMATEPRIORS}(\mathbf{Tr}, C)$ : { $C$  prior membership probability estimates}
8  $\mathbf{F} \leftarrow \emptyset$ 
9 parfor  $i \leftarrow 1$  to  $n$ 
10    $\mathbf{D}_i \leftarrow \text{BALANCEDBOOTSTRAPSAMPLE}(\mathbf{Tr})$ 
11   let  $\mathbf{D}_i = \langle \mathbf{P}, \mathbf{N}, \mathbf{U} \rangle$ 
12    $m_i \leftarrow \text{COMPUTE BBA}(C, \langle \mathbf{P}, \mathbf{N}, \mathbf{U} \rangle)$ 
13    $T_i \leftarrow \text{INDUCEETDTREE}(\mathbf{D}_i, C, \top, m, \hat{\mathbf{P}}_{\mathbf{r}})$ ;
14    $\mathbf{F} \leftarrow \mathbf{F} \cup \{T_i\}$ 
15 return  $\mathbf{F}$ 
16 end
```

427 *4.3.1. Growing ETRFs*

428 Alg. 3 describes the procedure for producing an ETRF. To this purpose, the
429 target concept C , a training set $\mathbf{Tr} \subseteq \text{Ind}(\mathcal{A})$ and the desired number of trees n
430 are required. \mathbf{Tr} may contain not only positive and negative examples but also
431 instances with uncertain membership w.r.t. C .

432 Similarly to a bagging approach, the training individuals are sampled with
433 replacement in order to obtain n subsets $\mathbf{D}_i \subseteq \mathbf{Tr}$, with $i = 1, \dots, n$. It is pos-
434 sible to apply various sampling strategies to obtain the various samples \mathbf{D}_i . In
435 this study we followed the same approach already used in our previous work [10].
436 Firstly, the initial data distribution is considered by adopting a stratified sam-
437 pling strategy w.r.t. the class-membership values to ensure the availability of
438 instances of the minority class. In the second phase, undersampling can be
439 performed on the training set in order to obtain (quasi-)balanced \mathbf{D}_i sets (i.e.
440 with a class imbalance that will not affect much the training process). This
441 means that if the majority class is the negative one, the exceeding part of the
442 negative examples is randomly discarded. In the dual case, positive instances
443 are removed. In addition, the sampling procedure removes also all the instances
444 of uncertain membership.

445 In Alg. 3, procedure BALANCEDBOOTSTRAPSAMPLE implements this strat-
446 egy returning the samples \mathbf{D}_i . For each \mathbf{D}_i , an ETDT T is built by invoking
447 the procedure INDUCEETDT. Note that, the procedure for learning an ETDT

Algorithm 4 Class-membership prediction routines for ETRFs

```

1 const  $\varepsilon \in ]0, 1]$  {decision threshold parameter}
2
3 function CLASSIFYBYETRF( $a, F, C$ ) :  $l$ 
4 input  $a$ : individual;  $F$ : ETRF;  $C$ : target concept
5 output  $l \in \mathcal{L}$ 
6 begin
7  $M[] \leftarrow$  new map {trees to BBAs}
8 parfor each  $T \in F$  do
9    $M[T] \leftarrow$  GETTREEBBA( $a, T$ )
10  $\bar{m} \leftarrow \bigoplus_{m \in M} m$  {pooling according to a combination rule}
11 for each  $\emptyset \neq s \in 2^\Omega$  do
12   Compute  $\overline{Bel}(s)$  from  $\bar{m}$ 
13 if  $|\overline{Bel}(\{-1\}) - \overline{Bel}(\{+1\})| \leq \varepsilon$  then
14   predlabel  $\leftarrow 0$  {case of uncertain membership}
15 else
16   predlabel  $\leftarrow \arg \max_{l \in \Omega} \overline{Bel}(\{l\})$  {cases of definite membership}
17 return predlabel
18 end
19
20 function GETTREEBBA( $a, T$ ) :  $\bar{m}$ 
21 input  $a$ : individual;  $T$ : ETDT
22 output  $\bar{m}$ : BBA
23 begin
24  $M \leftarrow$  GETLEAFBBALIST( $a, T$ ) {list of BBAs}
25  $\bar{m} \leftarrow \bigoplus_{m \in M} m$ 
26 return  $\bar{m}$ 
27 end

```

448 for the forest requires the introduction of some further amount of *randomiza-*
449 *tion*: the recursive case of Alg. 1 was modified so that the computation of the
450 BBAs and the selection of the best refinement are made considering a subset
451 $\mathbf{RS} \subseteq \rho(D)$ of randomly selected candidate specializations. This may be crucial
452 to improve the performance w.r.t. the one of a single classifier through a good
453 diversification among the trees.

454 4.3.2. Prediction

455 Given an ETRF, predictions can be made relying on the resulting classifica-
456 tion model. The related procedure, sketched in Alg. 4, works as follows.

457 Given the individual to be classified, for each tree T_i of the forest, the pro-
458 cedure GETTREEBBA returns a BBA obtained by pooling the various BBAs
459 found at the leaves reached from the root in a traversal path down the tree.

460 After polling all the trees in the ensemble, a set of BBAs deriving from the
461 previous phase are exploited to decide the classification for the test individual a .

462 Function CLASSIFYBYETRF takes an individual a and a forest F . Then, the
463 algorithm iterates on the forest trees collecting the BBAs via function GET-
464 TREEBBA.

465 Then, the BBAs are pooled according to a further combination rule, which
466 can be different from the one employed during the exploration of a single ETDT.
467 Additionally, this combination rule should be also an associative operator [19].
468 In this way, the result should not be affected by the pooling order of the BBAs.

469 In [18] we combined these BBAs via Dempster’s rule. Using this rule, the
470 disagreement among the classifiers, that corresponds to the conflict exploited
471 as a normalization factor, is explicitly considered by the meta-learner. Again,
472 the final decision is then made according to the belief function value computed
473 from the pooled BBAs \bar{m} .

474 4.4. Simplifying the Ensemble

475 In the previous works [10, 18], we noticed that a limited number of ETDTs
476 was usually sufficient to obtain a good performance. Growing forests with larger
477 numbers of trees did not improve significantly on predictiveness (in some cases
478 the performance even worsened). Moreover, the efficiency of the induction and
479 classification procedures obviously decayed owing to the increased number of
480 trees. Therefore, in this section, we illustrate how DST constructs can support
481 the simplification of an ETRF to increase the efficiency of the classification
482 phase while preserving its effectiveness.

483 The proposed solution (see Alg. 5) assumes that the prediction made using
484 an ETRF of progressively increasing size may lead to a poorer (or similar) per-
485 formance depending on the amount of conflictual evidence coming from a larger
486 number of trees. This basically implies that the confidence in the predictions
487 may decrease up to some point when the resulting predictions may even differ
488 from the expected ones.

489 The algorithm for pruning the ensemble is incremental, this means that it
490 works by considering one tree at a time. Specifically, given a forest \mathbf{F} , the
491 algorithm produces a new forest \mathbf{F}' as follows: it combines the pooled BBAs

Algorithm 5 Conflict-based ensemble simplification

```
1 function SIMPLIFICATION(F) : F'
2 input F: TRF
3 output F': TRF
4 begin
5  $M[] \leftarrow$  new array
6 for each  $T \in \mathbf{F}$  do
7      $M[T] \leftarrow$  GETBBAFROMTREE( $T$ )
8
9  $\bar{m} \leftarrow M[T_1]$ 
10 F'  $\leftarrow \{\}$  {initialize with the first ETDT in the forest}
11 for each  $T \in \mathbf{F}$  do
12      $c \leftarrow \sum_{B \cap C = \emptyset} \bar{m}(B)(M[T])(C)$ 
13     if  $c \leq \nu$  then
14         F'  $\leftarrow \mathbf{F}' \cup \{T\}$ 
15          $\bar{m} \leftarrow \bar{m} \oplus M[T]$ 
16
17 if F' =  $\emptyset$  then
18     F'  $\leftarrow \mathbf{F}' \cup \{T_1, T_2\}$  {return a forest with size=2, in case of oversimplification}
19
20 return F'
21 end
```

492 coming from each ETDT in the forest in order to compute the conflict measure c
493 (see Eq. 3). If the conflict does not go beyond a given threshold, namely ν , the
494 current tree T is added to **F'**.

495 The BBA drawn from a $T \in \mathbf{F}'$ and returned to the main procedure is
496 computed as follows: T is traversed following all the possible paths until all the
497 leaves are reached in order to collect the BBAs. Subsequently, the BBAs are
498 combined according to an associative rule (to avoid order-dependent results).
499 This is implemented in the procedure GETBBAFROMTREE. The resulting BBA
500 is then returned to the main procedure and used to determine c .

501 A particular case occurs when the conflict exceeds the threshold ν . In this
502 case, to prevent the production of an empty ETRF, the algorithm returns a
503 default forest composed by only two ETDTs.

504 5. Empirical Evaluation

505 The evaluation reported in this section aimed at assessing the effectiveness
506 of ETRFs and ETDTs proposed in this paper⁵.

⁵The source code is available at: <https://github.com/Giuseppe-Rizzo/SWMLAlgorithms>

Table 1: Ontologies employed in the experiments

Ontology	DL Lang.	#Axioms	#Concepts	#Roles	#Individuals
BCO	$\mathcal{ALCHO}\mathcal{F}(\mathcal{D})$	1098	196	22	112
BioPAX	$\mathcal{ALCL}\mathcal{F}(D)$	2617	74	70	323
NTN	$\mathcal{SH}\mathcal{I}\mathcal{F}(D)$	1516	47	27	676
HD	$\mathcal{ALCL}\mathcal{F}(D)$	8811	1498	10	639
FINANCIAL	$\mathcal{ALCL}\mathcal{F}(D)$	3509	60	16	1000
MONETARY	$\mathcal{ALCL}\mathcal{F}(D)$	7562	323	247	2466
DBPEDIA	\mathcal{ALCH}	78663	251	132	16606

Table 2: Distribution of the positive, negative and uncertain instances w.r.t the artificially generated target concepts for the various ontologies considered in the experiments

Ontology	% Pos.	% Neg.	% Unc.
BCO	17	53	30
BiOPAX	40	40	20
NTN	24	13	63
HD	24	11	65
FINANCIAL	26	47	30
MONETARY	36	44	20
DBPEDIA	16	14	70

507 5.1. Setup of the Experimental Sessions

508 The experiments have been carried out on various Web ontologies (see
509 Tab. 1) that are available on public repositories⁶. For each ontology, 15 query
510 concepts have been randomly generated by combining 2 through 8 (primitive
511 or defined) concepts of the ontology (using the conjunction and disjunction op-
512 erators or universal and existential restrictions). Each concept was generated
513 so that at least 40 positive examples and 40 negative examples can be found
514 among the individuals of the knowledge base.

515 Tab. 2 illustrates the average rate of the positive, negative, uncertain exam-
516 ples (computed considering all the individuals of $\text{Ind}(\mathcal{A})$) over the number of
517 query concepts.

⁶See <http://owl.cs.manchester.ac.uk/tools/repositories/>

518 We compared the methods and models proposed in this paper with a variety
519 of other approaches in the literature related to the task of inductive classification
520 with DL knowledge bases. Specifically, we selected:

- 521 • purely logical approaches, such as TDTs [3], CELOE [27], TRFs [10] and
522 the previous versions of ETDTs [7] and ETRFs;
- 523 • an instance-based method, i.e. the k -nearest neighbor algorithm embed-
524 ding a suitable distance measure as illustrated in [1];
- 525 • a kernel method for linear models, i.e. the *kernel perceptron* [28] adopting
526 a kernel function for individuals in DL knowledge bases [29, 1].

527 In the experiments with the ETDTs the three total uncertainty measures
528 reported in Sect. 3.2 have been considered: *non-specificity*, *confusion* and *dis-*
529 *sonance*. We repeated the experiments varying also the combination rules for
530 pooling the BBAs collected after tests with uncertain results are performed.
531 The rules adopted in the evaluation were: Dempster’s rule, *Dubois-Prade’s* rule
532 and the *mixing* rule.

533 The experiments on TRFs and ETRFs required a setup of the stratified
534 sampling rate and the forest size. Three sampling rates have been picked, 50%,
535 70% and 80%, while the forest size has been set to 10, 20 and 30 trees. In
536 the induction of (E)TDTs, the number of randomly selected specializations was
537 determined as the square root of candidate refinements: $n(C) = \sqrt{|\rho(C)|}$.
538 We ran the ETRF learning algorithm by varying three further parameters: the
539 heuristics for inducing the ETDTs, the combination rule for pooling the BBA
540 collected during the traversing process and the combination rules adopted as
541 meta-learner. Besides, we performed experiments with the ETRFs induction
542 algorithm with and without the simplification strategy, setting the threshold ν
543 to 0.4. Also, we set the value of parameter ε (Alg. 2 and Alg. 4) to 0.3 for forcing
544 the answer in favor of a definite membership, and the value for parameter θ
545 (Alg. 1), used to control the growth of a tree (either a TDT or an ETDT), was

546 heuristically⁷ set to 0.9.

547 Concerning the k -NN algorithm, we set the neighborhood size to $k = \log |\text{Tr}|$.
548 The distance measure between individuals have been chosen from the family
549 of measures proposed in [1] by setting its parameter p to 2 and using atomic
550 concepts in the signature of the knowledge base as a feature set.

551 In the experiments with CELOE, we set a *noise rate* of 25% (representing
552 the maximum number of admissible false negative cases).

553 Finally, the kernel perceptron required the choice of the *kernel function*, of
554 the *learning rate* and of the *number of epochs* for the training phase. In the
555 experiments, we used the kernel function between individuals of a DL knowledge
556 base proposed in [29, 1], and we set a learning rate of 0.05 and a number of
557 epochs of 200.

558 For each learning problem (each target concept considered for each
559 dataset/ontology), we estimated the average performance of the models under
560 comparison through a 10-fold cross validation procedure. The baseline (correct
561 classification labels) for the various instances in the training and test sets w.r.t.
562 the target concepts was computed by a DL reasoner. Specifically, the *macro-*
563 *averaged* F_1 -measure has been computed over the three membership values. In
564 addition, the following indices have been measured [3, 10, 7].

- 565 • match rate (M%), the percentage of test individuals for which the induc-
566 tive model agrees with the baseline (both positive, negative, or unknown);
- 567 • commission rate (C%), the fraction of test cases where the predicted mem-
568 bership is opposite w.r.t. the baseline (i.e. positive vs. negative or vice-
569 versa);
- 570 • omission rate (O%), the proportion of test cases for which the inductive
571 method cannot determine the definite membership that holds in the base-
572 line (i.e. unknown vs. positive or negative);

⁷ For each learning algorithm considered in the evaluation, the values have been tuned using a *leave-one-out* procedure.

- induction rate (I%). the percentage of test cases where the inductive method can predict a definite membership while it could not be determined for the baseline (i.e. positive or negative vs. unknown).

5.2. Outcomes

Table 3: Outcomes for ETDTs adopting the *mixing rule* in the classification step. The outcomes do not change significantly employing other combination rules

Ontology		NON-SPECIFICITY	DISSONANCE	CONFUSION
BCO	F_1	83.56 ± 05.06	84.15 ± 06.14	84.15 ± 06.14
	M%	85.48 ± 11.01	91.31 ± 14.79	91.31 ± 14.79
	C%	07.56 ± 08.08	00.86 ± 02.61	00.86 ± 02.61
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	I%	06.96 ± 05.97	07.83 ± 15.35	07.83 ± 15.35
BioPAX	F_1	82.16 ± 08.32	82.43 ± 06.47	86.98 ± 08.32
	M%	86.63 ± 14.60	87.00 ± 07.15	87.00 ± 07.15
	C%	11.02 ± 12.95	11.57 ± 02.62	11.57 ± 02.62
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	I%	02.35 ± 05.23	01.43 ± 08.32	01.43 ± 08.32
NTN	F_1	23.06 ± 26.14	14.65 ± 05.43	12.87 ± 26.54
	M%	23.87 ± 26.18	14.87 ± 24.18	13.85 ± 26.18
	C%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	I%	75.13 ± 26.18	85.13 ± 24.18	86.15 ± 26.17
HD	F_1	85.48 ± 11.01	91.31 ± 14.79	91.31 ± 14.79
	M%	10.69 ± 01.47	10.69 ± 01.47	10.69 ± 01.47
	C%	00.07 ± 00.17	00.07 ± 00.17	00.07 ± 00.17
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	I%	89.24 ± 01.46	89.24 ± 01.46	89.24 ± 01.46

Tables 3 through 17 present the outcomes of the various experiments. Preliminarily, note that for brevity, in the case of (E)TRFs, we report only the outcomes for ensemble models composed by 20 trees and induced using a 50% sampling rate as the performance had no significant variation in the experiments with the other values of such parameters. A similar consideration applies also to the experiments with (E)TDTs.

The results seem to be promising: ETDTs and ETRFs were competitive against other learning systems (see Tab. 3–4, 6–7, 8–9, 15–16). In some cases,

Table 4: Outcomes for ETDTs adopting the *mixing rule* in the classification step.

Ontology		NON-SPECIFICITY	DISSONANCE	CONFUSION
FINANCIAL	F_1	87.42 ± 08.23	88.23 ± 08.43	88.23 ± 08.43
	M%	83.43 ± 04.43	87.43 ± 17.42	87.43 ± 07.42
	C%	04.00 ± 03.35	00.00 ± 00.00	00.00 ± 00.00
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	I%	12.57 ± 13.45	07.53 ± 12.24	07.54 ± 12.24
MONETARY	F_1	85.48 ± 11.01	91.31 ± 14.79	91.31 ± 14.79
	M%	87.43 ± 13.45	93.47 ± 12.24	93.46 ± 12.24
	C%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	I%	12.57 ± 13.45	07.53 ± 12.24	07.54 ± 12.24
DBPEDIA	F_1	60.78 ± 23.08	60.78 ± 23.08	60.78 ± 23.08
	M%	53.84 ± 23.16	53.84 ± 23.16	54.46 ± 23.16
	C%	35.28 ± 23.30	35.28 ± 23.30	35.28 ± 23.30
	O%	00.00 ± 00.00	00.00 ± 00.00	00.54 ± 00.03
	I%	10.86 ± 01.69	10.86 ± 01.69	10.72 ± 13.30

585 the new models outperformed the others in terms of match rate, especially
586 K-NN and PERCEPTON. In general, we noted that the match rate obtained
587 with ETDTs was particularly high for ontologies endowed with a large number
588 of disjointness axioms, such as BIOPAX ,MONETARY and, to some extent, FI-
589 NANCIAL. The maximal average match rate were over 80% (for ETDTs) and
590 over 90% (for ETRFs). As regards the F_1 , it was particularly large on the
591 aforementioned ontologies and improved in the experiments with ETRFs.

Table 5: Average run-time (secs) for the classification using ETDTs varying the combination rule

Ontology	DEMPSTER	DUBOIS-PRADE	MIXING
BCO	3.5	3.5	2.4
BioPAX	7.5	7.5	7.5
NTN	4.5	4.5	2.5
HD	3.5	3.6	3
FINANCIAL	5	5	4.5
MONETARY	10.3	10.3	10.3
DBPEDIA	10	10	4

592 5.2.1. ETDTs

593 In the experiments with the ETDTs (see Tab. 3), the match rate was larger
594 and the commission rate was smaller using either the confusion or the dissonance
595 measure with respect to the outcomes observed when the non-specificity measure
596 was adopted. This was likely due to the fact that, adopting the non-specificity
597 measure, the heuristic basically tended to select concepts with a definite mem-
598 bership w.r.t. the target concept, with little or no increase of homogeneity in
599 the child nodes. As a consequence, even extending the branches with more de-
600 scendants, no significant gain was observed and the resulting ETDTs tended
601 to overfit the training instances. Conversely, adopting the confusion and dis-
602 sonance, the algorithm was biased towards shorter (and more predictive) trees
603 where pure leaves were obtained more easily.

604 As regards the employment of different combination rules to classify individ-
605 uals through ETDTs, we observed that none led to significant improvements.
606 On one hand, in the case of ontologies with properly defined constraints such
607 as concept disjointness, e.g. BIOPAX, the classification procedure tended to tra-
608 verse single branches thanks to intermediate tests with definite decisions. This
609 low degree of uncertainty yielded an analogous behavior w.r.t. the case of TDTs
610 and, consequently, to similar outcomes. On the other hand, in the case of
611 ontologies with a limited number of disjointness axioms more individuals ex-
612 hibited an uncertain membership w.r.t. the test concepts, so the classification
613 algorithm tends to traverse more branches reaching a larger number of BBAs
614 (at the leaves): the pooled BBAs obtained through the three combination rules
615 were very similar. Consequently, also the measures of belief used to decide
616 the final classification did not change significantly. This suggested that quasi-
617 associative rules, such as *mixing*, could be taken into account as alternative
618 strategies for combining evidence (despite their being order-dependent) that
619 are able to preserve the predictiveness of the classification models. This is an
620 advantage because classification through ETDTs via *mixing rule* was more ef-
621 ficient than with the adoption of the other rules. This benefit was particularly

Table 6: Outcomes for the ETRFs obtained adopting the three heuristics for the best concept selection and Dempster’s rule as meta-learner, with and without the use of the pruning

Ontology		No simplification			simplification		
		NON-SPECIFIC.	DISSONANCE	CONFUSION	NON-SPECIFIC.	DISSONANCE	CONFUSION
BCO	F_1	90.76 ± 06.67	91.76 ± 06.87	91.87 ± 07.23	95.23 ± 02.27	95.23 ± 02.27	95.23 ± 02.27
	M%	87.43 ± 09.13	88.23 ± 08.56	88.42 ± 08.43	92.31 ± 04.27	92.32 ± 04.27	92.31 ± 04.27
	C%	03.16 ± 03.09	02.44 ± 03.39	02.27 ± 03.38	02.81 ± 02.45	02.81 ± 02.45	02.91 ± 02.45
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	I%	09.41 ± 03.56	09.33 ± 03.46	09.31 ± 03.43	04.88 ± 03.45	04.88 ± 03.45	04.88 ± 03.45
BioPAX	F_1	90.37 ± 05.56	91.38 ± 05.57	92.43 ± 05.89	92.78 ± 05.84	92.78 ± 05.86	93.45 ± 04.57
	M%	93.45 ± 07.15	94.45 ± 07.14	94.45 ± 07.15	96.57 ± 06.15	95.98 ± 06.14	96.87 ± 06.23
	C%	05.22 ± 07.42	04.22 ± 07.42	04.22 ± 07.24	01.07 ± 01.67	01.71 ± 02.50	00.77 ± 01.74
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	I%	01.33 ± 07.16	01.33 ± 07.16	01.97 ± 07.16	02.36 ± 04.24	02.31 ± 04.13	02.30 ± 08.15
NTN	F_1	04.15 ± 03.25	04.15 ± 03.25	04.15 ± 03.25	35.25 ± 03.87	35.23 ± 03.87	35.23 ± 03.87
	M%	05.50 ± 07.28	05.50 ± 07.28	05.50 ± 07.28	26.40 ± 05.15	26.43 ± 05.15	26.43 ± 05.14
	C%	06.52 ± 07.54	06.52 ± 07.54	06.52 ± 07.54	06.52 ± 07.54	06.52 ± 07.54	06.52 ± 07.54
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	I%	87.99 ± 08.84	87.99 ± 08.84	87.99 ± 08.84	67.05 ± 05.35	67.05 ± 05.35	67.07 ± 05.35
HD	F_1	08.32 ± 00.15	08.32 ± 00.15	08.32 ± 00.15	28.15 ± 02.17	28.15 ± 02.17	28.15 ± 02.17
	M%	10.29 ± 00.00	10.29 ± 00.01	10.29 ± 00.02	32.56 ± 00.43	33.43 ± 00.43	33.56 ± 00.42
	C%	00.57 ± 00.05	00.57 ± 00.05	00.57 ± 00.05	00.14 ± 00.26	00.14 ± 00.27	00.14 ± 00.28
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	I%	89.14 ± 00.26	89.14 ± 00.26	89.14 ± 00.26	67.44 ± 00.26	67.44 ± 00.26	67.44 ± 00.26

622 evident in the experiments with larger ontologies like DBPEDIA. For this ontol-
623 ogy, on average the classification of an individual adopting the mixing rule was
624 60% faster than with the Dubois-Prade rule. Tab. 5 summarizes the average
625 time (in seconds) required by an ETDT for classifying an individual.

626 5.2.2. ETRFs and simplification procedure

627 Concerning the experiments with the ETRFs (Tab. 6 - 8), as expected, the
628 ensemble models showed on average a superior performance with respect to the
629 ETDTs, and in most of the cases we observed a decrease of standard deviation.
630 As previously mentioned, the F_1 increased with respect to the experiments with
631 ETDTs, suggesting that the sampling strategy brought benefits to the predic-
632 tiveness of ETRFs mitigating the of bias classification models towards the most

Table 7: Outcomes for the ETRFs obtained adopting the three heuristics for the best concept selection and Dempster’s rule as meta-learner, with and without the use of the pruning

Ontology		No simplification			simplification		
		NON-SPECIFIC.	DISSONANCE	CONFUSION	NON-SPECIFIC.	DISSONANCE	CONFUSION
FINANCIAL	F_1	90.14 ± 06.76	92.46 ± 07.43	96.79 ± 03.17	96.79 ± 03.17	96.79 ± 03.17	96.79 ± 03.17
	M%	93.43 ± 05.06	93.89 ± 05.16	94.03 ± 05.23	97.12 ± 03.10	97.13 ± 04.12	97.12 ± 04.15
	C%	01.07 ± 01.67	01.71 ± 02.50	00.77 ± 01.74	00.60 ± 00.03	00.54 ± 00.03	00.77 ± 01.74
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	I%	02.28 ± 08.13	02.31 ± 08.17	02.30 ± 08.15	02.28 ± 08.13	02.31 ± 08.17	02.30 ± 08.15
MONETARY	F_1	95.23 ± 03.24	96.45 ± 03.76	96.57 ± 03.17	97.43 ± 02.14	97.43 ± 02.14	97.43 ± 02.14
	M%	93.43 ± 05.06	95.89 ± 05.16	94.56 ± 04.46	96.65 ± 04.35	99.43 ± 08.13	99.55 ± 08.15
	C%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	I%	06.57 ± 05.06	04.11 ± 05.16	01.97 ± 07.16	03.35 ± 04.35	00.57 ± 08.13	00.45 ± 08.15
DBPEDIA	F_1	60.43 ± 02.15	60.43 ± 02.15	61.25 ± 02.24	68.34 ± 02.15	68.43 ± 02.15	68.34 ± 02.15
	M%	53.84 ± 05.43	53.84 ± 05.43	54.46 ± 05.43	70.44 ± 03.31	70.43 ± 03.31	70.44 ± 03.31
	C%	00.08 ± 00.01	00.08 ± 00.02	00.08 ± 00.01	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	I%	45.28 ± 23.30	45.28 ± 23.30	45.28 ± 23.30	29.56 ± 03.31	29.57 ± 03.31	29.56 ± 03.31

633 probable membership value. Such a stability of the ensemble models was likely
634 due to the mediation operated by the the meta-learner over the various models,
635 that positively influenced the final decision towards correct label assignments.
636 Again, the choice of the combination rule for the BBAs at the leaves of a sin-
637 gular ETDT did not affect significantly the performance of the considered tree
638 models. Conversely, using further combination rules as meta-learners had a
639 stronger influence on the performance. In particular, adopting Dubois-Prade
640 rule we observed a decrease of the induction rate and an increase of the match
641 rate. A similar outcome was obtained using the mixing rule. Unlike Dempster’s
642 rule, the adoption of the Dubois-Prade and mixing rules tended to reduce the
643 evidence in favor of a definite membership. This means that the belief related
644 to the hypotheses of positive and negative memberships were generally low and
645 their difference often did not exceed the threshold ϵ .

646 A similar effect was observed in the experiments with the simplifica-
647 tion method proposed in the paper: smaller ensembles tended to predict an

Table 8: Outcomes for ETRFs adopting the three heuristics for the best concept selection and the Dubois-Prade rule as a meta-learner, with and without the use of the pruning

Ontology	No simplification			simplification			
	NON-SPECIFIC.	DISSONANCE	CONFUSION	NON-SPECIFIC.	DISSONANCE	CONFUSION	
BCO	F_1	87.13 ± 05.67	90.45 ± 03.56	90.48 ± 03.78	89.94 ± 02.13	89.13 ± 02.19	89.13 ± 02.15
	M%	90.44 ± 09.13	93.24 ± 08.56	93.40 ± 08.35	93.17 ± 04.27	94.45 ± 04.27	94.44 ± 04.15
	C%	03.16 ± 03.09	02.43 ± 03.39	02.29 ± 03.45	02.81 ± 02.45	02.81 ± 02.45	02.91 ± 02.45
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	I%	06.40 ± 03.56	04.33 ± 03.27	04.31 ± 03.46	04.01 ± 03.45	02.73 ± 03.45	02.74 ± 03.56
BioPAX	F_1	90.98 ± 03.79	92.45 ± 03.79	92.45 ± 03.79	93.76 ± 04.25	94.33 ± 05.25	94.33 ± 05.25
	M%	93.45 ± 07.15	94.45 ± 07.14	94.45 ± 07.15	96.57 ± 06.15	95.98 ± 06.14	96.87 ± 06.23
	C%	05.22 ± 07.42	04.22 ± 07.42	04.22 ± 07.24	01.07 ± 01.67	01.71 ± 02.50	00.77 ± 01.74
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	I%	01.33 ± 07.16	01.33 ± 07.16	01.97 ± 07.16	02.36 ± 04.24	02.31 ± 04.13	02.36 ± 08.15
NTN	F_1	47.98 ± 03.46	47.98 ± 03.46	47.98 ± 03.46	56.78 ± 03.24	56.78 ± 03.24	56.78 ± 03.24
	M%	57.68 ± 03.43	57.68 ± 03.43	57.68 ± 03.43	60.40 ± 05.45	60.40 ± 05.45	60.40 ± 05.45
	C%	06.52 ± 07.54	06.52 ± 07.54	06.55 ± 07.54	06.55 ± 07.54	06.55 ± 07.55	06.55 ± 07.55
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	I%	35.88 ± 08.84	35.88 ± 08.84	35.88 ± 08.84	33.05 ± 05.73	33.05 ± 05.73	33.05 ± 05.73
HD	F_1	50.44 ± 00.13	50.44 ± 00.13	50.44 ± 00.13	65.43 ± 00.35	67.43 ± 00.43	65.43 ± 00.13
	M%	59.49 ± 00.03	59.49 ± 00.03	59.49 ± 00.03	67.56 ± 00.43	68.43 ± 00.43	67.56 ± 00.42
	C%	00.47 ± 00.05	00.47 ± 00.05	00.47 ± 00.05	00.14 ± 00.26	00.14 ± 00.27	00.14 ± 00.28
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	I%	40.04 ± 00.26	40.04 ± 00.26	40.04 ± 00.26	32.30 ± 00.26	32.30 ± 00.26	32.30 ± 00.26

648 uncertain-membership more easily than the forests obtained without the appli-
649 cation of the pruning strategy. However, thanks to simplification strategy, the
650 size of the resulting forests was considerably reduced: after pruning, the average
651 size of the ETRFs did not exceed 10 trees. Tab. 12 reports the average forest
652 sizes⁸.

653 5.2.3. Evaluating cases of induction

654 One of the most important consequences of the *credulous* behavior of ETDTs
655 and ETRFs was the large induction rates, which represent the cases of non

⁸The sizes have been averaged over the folds and, the resulting values have been further averaged over the number of target concepts.

Table 9: Outcomes for ETRFs adopting the three heuristics for the best concept selection and the Dubois-Prade rule as a meta-learner, with and without the use of the pruning

Ontology	No simplification			simplification			
	NON-SPECIFIC.	DISSONANCE	CONFUSION	NON-SPECIFIC.	DISSONANCE	CONFUSION	
FINANCIAL	F_1	90.89 ± 03.25	90.97 ± 03.36	91.23 ± 03.76	96.85 ± 03.25	96.85 ± 03.25	96.85 ± 03.25
	M%	93.43 ± 05.06	93.89 ± 05.16	94.03 ± 05.23	97.12 ± 03.10	97.13 ± 04.12	97.12 ± 04.15
	C%	01.07 ± 01.67	01.71 ± 02.50	00.77 ± 01.74	00.60 ± 00.03	00.54 ± 00.03	00.77 ± 01.74
	O%	03.22 ± 00.15	02.19 ± 00.15	02.90 ± 00.14	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	I%	02.28 ± 08.13	02.31 ± 08.17	02.30 ± 08.15	02.28 ± 08.13	02.31 ± 08.17	02.30 ± 08.15
MONETARY	F_1	92.39 ± 04.97	94.76 ± 05.76	94.45 ± 07.15	95.57 ± 06.15	97.21 ± 05.67	97.43 ± 03.35
	M%	93.43 ± 05.06	95.89 ± 05.16	94.56 ± 04.46	96.65 ± 04.35	99.43 ± 08.13	99.55 ± 08.15
	C%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	I%	06.57 ± 05.06	04.11 ± 05.16	01.97 ± 07.16	03.35 ± 04.35	00.57 ± 08.13	00.45 ± 08.15
DBPEDIA	F_1	59.89 ± 03.78	59.89 ± 03.78	50.23 ± 02.43	68.12 ± 03.24	68.12 ± 03.24	68.12 ± 03.24
	M%	63.84 ± 05.43	63.84 ± 05.43	54.46 ± 05.43	70.43 ± 03.31	70.43 ± 03.31	70.43 ± 03.31
	C%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	I%	35.16 ± 23.30	35.16 ± 23.30	35.16 ± 23.30	29.56 ± 03.31	29.56 ± 03.31	29.56 ± 03.31

656 logically derivable definite classifications whose correctness requires a validation
657 from a domain expert.

658 However, the number of new assertions resulting from the inductive classifi-
659 cation models was very large, especially in the experiments with NTN, HD, and
660 DBPEDIA. As a consequence, there was a drastic decrease of the F -measure as
661 it considers such cases as label mismatches, whereas the induction rate treats
662 them as non conflictual assertions that could be exploited in a perspective of
663 integration and evolution of the KBs.

664 Devising a different strategy for tackling these cases of induction, we de-
665 signed and performed new experiments, considering a modified version of the
666 ontologies. The new versions were obtained by introducing disjointness axioms
667 in accordance with the *strong disjointness assumption* (SDA) which states that
668 sibling concepts in the subsumption hierarchy can be considered as disjoint [30].
669 In this way, the cases of individuals with uncertain-membership can be mini-
670 mized or totally avoided and a ground truth with definite membership labels can

Table 10: Outcomes for ETDTs under *Strong Disjointness Assumption*

Ontology		NON-SPECIF.	DISSONANCE	CONFUSION
NTN	F_1	92.17 ± 07.56	93.78 ± 07.43	93.78 ± 07.43
	M%	93.45 ± 07.67	94.67 ± 07.85	94.67 ± 07.85
	C%	06.55 ± 07.67	05.33 ± 07.85	05.33 ± 07.85
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
HD	F_1	94.12 ± 03.57	94.12 ± 03.57	94.12 ± 03.57
	M%	96.46 ± 04.56	96.46 ± 04.56	96.46 ± 04.56
	C%	03.54 ± 04.56	03.54 ± 04.56	03.54 ± 04.56
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
DBPEDIA	F_1	91.05 ± 02.43	91.05 ± 02.43	91.05 ± 02.43
	M%	92.35 ± 03.97	92.35 ± 03.97	92.35 ± 03.97
	C%	07.65 ± 03.97	07.65 ± 03.97	07.65 ± 03.97
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00

671 be provided to evaluate induction cases. Tab. 10 and 11 illustrate the results
672 of the new experiments with NTN, HD and DBPEDIA. Note that, under the
673 SDA, most of the cases previously classified as cases of induction were deemed
674 as matching cases via both ETDTs and ETRFs. Again, the performance of the
675 ETRFs overcame the one obtained through a single tree in terms of F_1 , match
676 rate and also with a decrease of the standard deviation. With the adoption of
677 the SDA, the match rates were less biased by the parameter ϵ tuned for ETDTs
678 and ETRFs.

679 *5.2.4. Experiments with ETDTs and ETRFs with Special Probabilistic BBAs*

680 For the sake of completeness, we tested the effectiveness of a modified version
681 of the ETDT and ETRF models (and related algorithms) such that the BBAs in
682 their nodes had only singletons as focal elements. To this purpose, the function
683 COMPUTEBBA in Alg. 1 has been adapted as follows: the probability mass
684 assigned to $m(\Omega)$ has been proportionally distributed to the singletons $\{-1\}$ and
685 $\{+1\}$ (preserving the sum (1) of the focal elements).

686 Similarly to the previous experiments, Tab. 13 and 14 illustrate the outcomes
687 of this comparison only for NTN, HD and DBPEDIA ontologies, where the
688 results significantly changed w.r.t. the original versions (in the experiments with

Table 11: Outcomes for the ETRFs under *Strong Disjointness Assumption*

Ontology		No simplification			simplification		
		NON-SPECIF.	DISSONANCE	CONFUSION	NON-SPECIF.	DISSONANCE	CONFUSION
NTN	F_1	96.23 ± 03.13	96.32 ± 04.43	96.32 ± 04.18	95.87 ± 04.56	96.74 ± 03.85	96.74 ± 03.85
	M%	96.57 ± 04.23	96.60 ± 04.17	96.60 ± 04.17	95.87 ± 04.56	96.74 ± 03.85	96.74 ± 03.85
	C%	03.43 ± 04.23	03.40 ± 04.17	03.40 ± 04.17	04.13 ± 04.56	03.26 ± 03.85	03.26 ± 03.85
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
HD	F_1	97.23 ± 00.16	97.43 ± 00.16	97.43 ± 00.17	97.23 ± 00.16	97.43 ± 00.16	97.43 ± 00.16
	M%	98.56 ± 00.43	98.80 ± 00.45	98.70 ± 00.34	98.60 ± 00.44	98.76 ± 00.32	98.76 ± 00.33
	C%	01.44 ± 00.43	01.20 ± 00.45	01.30 ± 00.34	01.40 ± 00.44	01.24 ± 00.32	01.24 ± 00.33
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
DBPEDIA	F_1	99.43 ± 00.12	99.43 ± 00.12	99.43 ± 00.12	99.43 ± 00.12	99.43 ± 00.12	99.22 ± 03.14
	M%	99.20 ± 03.21	99.26 ± 03.17	99.16 ± 03.21	99.21 ± 03.13	99.21 ± 03.12	99.22 ± 03.14
	C%	00.80 ± 03.21	00.74 ± 03.17	00.84 ± 03.21	00.70 ± 03.13	00.79 ± 03.12	00.78 ± 03.14
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00

Table 12: Average size of forests (number of trees) after the pruning

Ontology	Forest size after the pruning		
	10 trees	20 trees	30 trees
BCO	6.76	6.43	5.76
BioPAX	5.56	5.43	5.76
NTN	7.87	7.45	7.43
HD	8.43	7.65	6.56
FINANCIAL	4.34	4.43	4.43
MONETARY	8.44	7.65	7.42
DBPEDIA	8.44	7.23	7.33

689 the other ontologies, the results did not change because the BBAs of the trees
690 have already singletons as focal elements). The tables report the outcomes
691 obtained inducing ETDTs and ETRFs with the mixing rule (for pooling the
692 BBAs in the leaf-nodes) and Dubois-Prade’s rule as a meta-learner. Similar
693 values have been obtained in the evaluation with the other rules.

694 Generally speaking, we noted a decay of the performance in terms of F -
695 measure, both for ETDTs and ETRFs w.r.t. the original versions, especially in
696 terms of (an increased) commission rate. On a close inspection of the models,
697 we observed that the BBAs at the leaves nodes computed with the new proce-

Table 13: Outcomes for ETDTs with BBAs having singletons as focal elements

Ontology		NON-SPECIF.	DISSONANCE	CONFUSION
NTN	F_1	73.23 ± 12.54	76.24 ± 12.43	73.15 ± 12.44
	M%	73.42 ± 11.43	77.32 ± 07.85	74.23 ± 07.85
	C%	13.32 ± 12.43	09.15 ± 04.34	09.15 ± 03.85
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	I%	13.26 ± 08.84	13.26 ± 08.84	13.26 ± 08.84
HD	F_1	67.16 ± 13.43	67.16 ± 13.43	67.16 ± 13.43
	M%	70.01 ± 07.26	70.01 ± 07.26	70.01 ± 07.26
	C%	14.65 ± 04.56	14.65 ± 04.56	14.65 ± 04.56
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	I%	15.33 ± 01.35	15.34 ± 01.35	15.34 ± 01.34
DBPEDIA	F_1	86.43 ± 03.35	86.43 ± 03.35	86.43 ± 03.35
	M%	87.85 ± 13.23	87.85 ± 13.23	87.85 ± 13.23
	C%	00.37 ± 00.30	00.37 ± 00.29	00.37 ± 00.29
	O%	00.30 ± 00.06	00.30 ± 00.06	00.30 ± 00.06
	I%	11.27 ± 08.73	11.27 ± 08.73	11.27 ± 08.73

698 dure tended to favor the majority class with high assigned values. When such
 699 functions are pooled through a combination rule, the final decision was strongly
 700 biased towards such class. As a consequence, the models determined a wrong
 701 membership value for the test individuals. Another remarkable difference is in
 702 the lower induction rate, that was likely due to the induction of trees in which
 703 the focal values of the BBAs, $m(\{+1\})$ and $m(\{-1\})$, located in the leaf-nodes
 704 were often (approximately) equal. Such cases represented the main source for
 705 ties, resulting in a label 0 returned. In this sense, even resorting to forests in-
 706 stead of single trees did not allow to considerably improve the performance: the
 707 membership assessed by one tree was further confirmed by the other trees in
 708 the forest.

709 *5.2.5. Comparison with other inductive systems*

710 As previously described, ETDTs and ETRFs showed a more credulous be-
 711 havior w.r.t. the other learning systems used in the experiments, in particular
 712 compared to the instance-based methods and CELOE (see Tab. 15 and 16).
 713 The k -NN showed a very *cautious* behavior: the neighborhood of the test in-

Table 14: Outcomes for the ETRFs with BBAs having only singletons as focal elements

Ontology	No simplification			simplification			
	NON-SPECIF.	DISSONANCE	CONFUSION	NON-SPECIF.	DISSONANCE	CONFUSION	
NTN	F_1	74.08 ± 08.15	74.08 ± 08.15	74.08 ± 08.15	75.16 ± 10.54	75.16 ± 10.54	75.16 ± 10.54
	M%	75.34 ± 09.23	75.23 ± 09.23	75.24 ± 09.24	76.87 ± 09.14	76.88 ± 09.14	76.88 ± 09.14
	C%	13.32 ± 12.21	13.43 ± 12.21	13.42 ± 12.20	13.32 ± 12.21	13.32 ± 12.21	13.32 ± 12.21
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	I%	11.34 ± 08.76	11.34 ± 08.76	11.34 ± 08.84	09.81 ± 04.23	09.81 ± 04.23	09.81 ± 04.23
HD	F_1	73.25 ± 07.42	73.26 ± 07.42	73.25 ± 07.42	73.25 ± 07.42	73.25 ± 07.42	73.25 ± 07.42
	M%	74.32 ± 05.13	74.32 ± 05.13	74.32 ± 05.13	74.32 ± 05.13	74.32 ± 05.13	74.32 ± 05.13
	C%	10.31 ± 04.56	10.31 ± 04.56	10.31 ± 04.56	10.31 ± 04.56	10.31 ± 04.56	10.31 ± 04.56
	O%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	I%	15.33 ± 01.35	15.34 ± 01.35	15.34 ± 01.34	15.33 ± 01.35	15.34 ± 01.35	15.34 ± 01.34
DBPEDIA	F_1	86.43 ± 03.35	86.43 ± 03.35	86.43 ± 03.35	86.43 ± 03.35	86.43 ± 03.35	86.43 ± 03.35
	M%	87.85 ± 13.23	87.85 ± 13.23	87.85 ± 13.23	87.85 ± 13.23	87.85 ± 13.23	87.85 ± 13.23
	C%	00.37 ± 00.30	00.37 ± 00.29	00.37 ± 00.29	00.37 ± 00.30	00.37 ± 00.29	00.37 ± 00.29
	O%	00.30 ± 00.06	00.30 ± 00.06	00.30 ± 00.06	00.30 ± 00.06	00.30 ± 00.06	00.30 ± 00.06
	I%	11.27 ± 08.73	11.27 ± 08.73	11.27 ± 08.73	11.27 ± 08.73	11.27 ± 08.73	11.27 ± 08.73

714 individuals was often made up of uncertain individuals. This explains both the
715 very high match rate achieved with this algorithm in the experiments with NTN
716 and the high omission rate observed in the experiments with DBPEDIA. In the
717 experiments with CELOE, introducing a stricter definition of negative example
718 than the one originally adopted in [27], made the algorithm more sensitive to
719 lack of disjointness axioms and, consequently, led to omission cases rather than
720 commission errors. Conversely, in case of ontologies with an explicit specifica-
721 tion of disjointness axioms, the match rate tended to be very high (in some cases
722 close to 100%), thanks to a strategy that aims at maximizing the F -measure.

723 Finally, in the experiments with PERCEPTRON, we observed a drop of the
724 match rate and an increase of commission and induction cases. On one hand, the
725 higher commission rates were due to overfitting models, likely owing to the large
726 number of epochs adopted in the experiments. On the other hand, the higher
727 induction rates were due to the decision procedure adopted in the classification
728 phase, which tended to assign a definite membership rather than an uncertain
729 membership to test individuals.

Table 15: Outcomes for other learning systems

Ontology		TDT	TRF	k-NN	CELOE	PERCEPTRON
BCO	F_1	76.23 ± 03.01	84.78 ± 02.43	84.78 ± 02.43	100.0 ± 00.00	83.45 ± 12.45
	M%	80.44 ± 11.01	87.99 ± 07.85	87.83 ± 12.43	100.0 ± 00.00	86.27 ± 15.79
	C%	07.56 ± 08.08	04.32 ± 04.68	12.77 ± 04.77	00.00 ± 00.00	02.47 ± 03.70
	O%	05.04 ± 04.28	00.09 ± 00.27	00.02 ± 00.04	00.00 ± 00.00	00.00 ± 00.00
	I%	06.96 ± 05.97	07.61 ± 06.82	00.40 ± 00.00	00.00 ± 00.00	09.36 ± 13.96
BioPAX	F_1	64.23 ± 13.26	71.43 ± 03.24	77.23 ± 03.46	100.0 ± 00.00	63.43 ± 15.46
	M%	66.63 ± 14.60	75.93 ± 17.05	75.49 ± 17.05	75.30 ± 16.23	65.30 ± 16.23
	C%	31.03 ± 12.95	22.11 ± 16.54	18.54 ± 17.80	18.74 ± 17.80	18.74 ± 17.80
	O%	00.39 ± 00.61	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	I%	01.95 ± 07.13	01.97 ± 07.16	01.97 ± 07.16	01.97 ± 07.16	11.97 ± 05.76
NTN	F_1	63.24 ± 10.98	81.43 ± 03.35	95.23 ± 03.23	78.91 ± 08.43	96.14 ± 08.43
	M%	68.85 ± 13.23	83.42 ± 07.85	96.82 ± 03.43	83.42 ± 07.85	96.81 ± 07.46
	C%	00.37 ± 00.30	00.00 ± 00.00	00.02 ± 00.04	00.02 ± 00.04	00.02 ± 00.04
	O%	09.51 ± 07.06	13.40 ± 10.17	00.00 ± 00.00	00.02 ± 00.04	00.00 ± 00.00
	I%	21.27 ± 08.73	03.16 ± 04.65	03.16 ± 04.65	00.00 ± 00.00	03.17 ± 04.65
HD	F_1	54.23 ± 14.15	62.85 ± 10.43	66.54 ± 17.11	65.00 ± 17.63	66.42 ± 16.43
	M%	58.31 ± 14.06	67.95 ± 16.99	67.96 ± 17.00	67.95 ± 17.03	68.00 ± 16.98
	C%	00.44 ± 00.47	00.02 ± 00.05	00.01 ± 00.05	00.02 ± 00.05	00.02 ± 00.05
	O%	05.51 ± 01.81	06.38 ± 02.03	06.38 ± 02.03	06.38 ± 02.03	06.38 ± 02.03
	I%	35.74 ± 15.90	25.61 ± 18.98	25.61 ± 18.98	25.61 ± 18.98	25.59 ± 18.98

730 5.2.6. Efficiency of the methods

731 A final remark is related to the efficiency of the proposed approaches. Con-
732 sidering Tab. 17 it can be noted that the averaged run-times of the ETDT and
733 ETRF models spanned from less than 35s to almost 13000s. The efficiency of
734 the solutions proposed in this paper depends on the size of training sets and
735 the number of concepts and roles contained in the signature of the knowledge
736 bases. While the former affected the performance in terms of the number of
737 tests to be performed in the training/test phase, which was intensively used by
738 ETDTs and ETRFs, the latter affected the generation of the complex concept
739 descriptions installed into the nodes. Also, the pruning procedure employed for
740 optimizing the ensemble models represented a further complexity source in the
741 training phase but simpler models brought an increased efficiency in the pre-
742 diction phase. Overall, the efficiency of the new models in both training and

Table 16: Outcomes for other learning systems

Ontology		TDT	TRF	K-NN	CELOE	PERCEPTRON
FINANCIAL	F_1	66.23 ± 36.01	96.23 ± 02.56	96.23 ± 02.56	99.12 ± 00.73	74.32 ± 00.87
	M%	67.06 ± 36.09	96.70 ± 00.48	96.70 ± 00.65	99.70 ± 00.68	79.50 ± 00.68
	C%	00.00 ± 00.00	02.00 ± 03.43	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	O%	32.94 ± 36.09	00.00 ± 00.60	00.30 ± 00.68	00.30 ± 00.68	00.00 ± 00.00
	I%	00.00 ± 00.00	01.30 ± 00.50	00.00 ± 00.00	00.00 ± 00.00	20.50 ± 00.68
MONETARY	F_1	66.12 ± 15.23	94.13 ± 07.74	100.0 ± 00.00	100.0 ± 00.00	65.43 ± 15.96
	M%	68.93 ± 15.87	94.53 ± 07.68	100.0 ± 00.00	100.0 ± 00.00	68.93 ± 15.87
	C%	06.14 ± 07.20	05.47 ± 07.68	00.00 ± 00.00	00.00 ± 00.00	06.14 ± 07.20
	O%	16.94 ± 09.74	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
	I%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00
DBPEDIA	F_1	08.12 ± 01.21	08.12 ± 01.21	08.12 ± 01.21	56.13 ± 20.43	58.12 ± 15.47
	M%	10.86 ± 01.69	10.86 ± 01.69	10.86 ± 01.69	58.84 ± 20.35	63.93 ± 15.07
	C%	43.12 ± 00.57	43.12 ± 00.57	43.12 ± 00.57	30.28 ± 20.10	25.18 ± 14.48
	O%	46.02 ± 01.64	46.02 ± 01.69	46.02 ± 01.69	00.00 ± 00.00	00.00 ± 00.00
	I%	00.00 ± 00.00	00.00 ± 00.00	00.00 ± 00.00	10.86 ± 01.69	10.86 ± 01.69

743 test phase is comparable to the one of TDTs and TRFs and also close to the
744 average execution time with the k -NN. Indeed, one of the main bottlenecks of
745 the lazy learning approach was related to the (exhaustive) search of the near-
746 est neighbors for each test individual. Moreover, we noted that the evidential
747 models were more efficient than PERCEPTRON. In this case, the run-times of
748 PERCEPTRON were affected mainly by the inefficiency of the training phase in
749 which, for each epoch, all the training examples are processed to determine the
750 coefficients of the classification model.

751 6. Related Work

752 The knowledge made available in a decentralized form across the Semantic
753 Web is often contradictory, imprecise and incomplete [31]. Machine learning can
754 be exploited for setting up methods providing alternative forms of reasoning.
755 In this work, we specifically focused on the task of assessing the membership of
756 an individual with respect to a target concept. This problem has been largely
757 investigated in the literature and various approximate classification models have

Table 17: Ranges of average run-time (training / test) per experiment (s)

Ontology	ETDT	ETRF	TDT	TRF	K-NN	CELOE	PERCEPTRON
	[min, max]	[min, max]	[min, max]	[min, max]	[min, max]	[min, max]	[min, max]
BCO	[35, 40]	[120, 453]	[35, 40]	[120, 452]	[65, 87]	[15, 37]	[480,520]
BioPAX	[67, 87]	[123, 456]	[70, 103]	[145, 523]	[83, 245]	[23, 60]	[345, 725]
NTN	[432, 785]	[876, 1256]	[578, 876]	[914, 1243]	[123, 456]	[12,65]	[765, 1343]
HD	[446, 879]	[1245, 1278]	[446, 895]	[1567, 1568]	[123, 456]	[12,65]	[1234, 1456]
FINANCIAL	[1845, 2567]	[18765,42345]	[1845, 2567]	[18765,42345]	[12456,12876]	[87,247]	[24569,56797]
MONETARY	[2476, 4587]	[3687,45890]	[2444, 4598]	[3687,45890]	[14876,15321]	[124,256]	[49872,58931]
DBPEDIA	[3211, 4237]	[12345,12789]	[3211, 4237]	[12345,12789]	[8769,14321]	[87,231]	[23456,60432]

758 been proposed [1].

759 *Non-parametric* methods are among the most common solutions. Among
760 the others, the *k*-nearest neighbor procedure [1] (also employed in the experi-
761 ments) and the *reduced Coulomb energy network* [32] have been proposed. Both
762 approaches exploit a language-independent distance measure between individ-
763 uals in a DL knowledge base. Such a metric is computed based on a set of
764 projection functions that express the behavior of an individual w.r.t. a set of
765 concepts (treated as logic *features*). Essentially the aim is selecting prototypi-
766 cal individuals and classifying unseen ones on the ground of the similarity w.r.t.
767 the closest prototypes, the neighborhood, that in the latter case is mediated
768 by a network model (similar to the *radial basis function networks* [33]). Other
769 related solutions are based on the explicit adaptation of *kernel methods*. For
770 instance, in the evaluation, we used the kernel *perceptron* [28] adopting a kernel
771 function that is closely related to the distance measure adopted by the classifiers
772 described above [29].

773 Other solutions stem from *concept learning* algorithms devised in ILP to
774 solve a closely related problem. The goal is to obtain an explicit intensional
775 definition (a concept description in terms of the language bias of choice) de-
776 scribing the available examples that should be general enough to account also
777 for unseen instances. Various algorithms of this kind have been proposed, e.g.
778 DL-FOIL [34], CELOE [27] and the mentioned method for the induction of
779 terminological decision trees [3]. [The latter extend decision trees for multi-](#)

780 relational representations (such as first-order logic fragments [4, 5] and *selection*
781 *graphs* [35, 36]) towards SW representations. A related approach, based on
782 models called *Semantic Decision Trees*, has been proposed in [37]. Although
783 they are indeed quite similar to the mentioned TDTs [3], their empirical evalu-
784 ation did not compare these models and it was limited to very small knowledge
785 bases. All these approaches are based on the use of a refinement operator in or-
786 der to progressively build such description(s). However, such concept learning
787 methods often do not provide a strategy for representing uncertainty, although
788 various efforts have been devoted to investigate the effectiveness of models com-
789 bining multi-relational representation languages and uncertainty, in the context
790 of *Statistical Relational Learning* [38] or *Probabilistic Inductive Logic Program-*
791 *ming* [39]. Among the existing solutions it is possible to mention *Bayesian*
792 *Logic Programs* [40] and *Markov Logic Networks* (MLNs) [41]. Focusing on
793 MLNs, a *domain closure* assumption is required thus diverting from the open-
794 world semantics of the FOL fragments adopted as standard representations in
795 the SW context [42]. However, the assumptions for inducing MLNs can be re-
796 laxed by using an EM algorithm to learn from incomplete data [43]. In this
797 perspective, a recent work [44] has proposed a functional-gradient boosting al-
798 gorithm based on EM in order to learn, under the OWA, the structure and the
799 parameters of the models simultaneously.

800 The need to circumvent the exponential growth of the model (and hence
801 of the number of parameters) required by the groundings justified the works
802 on *approximation methods* [45] and *lifted inference techniques* [46, 47] Alterna-
803 tively, *tensor models* have also been proposed [48, 49] although the limitations
804 in terms of scalability of such complex statistical models remains. That is why
805 currently *representation learning* approaches [50] have attracted the attention
806 of the community. They trade the focus on the mere relational structure of the
807 rich SW KBs with a low rank representation which is more manageable with
808 standard geometric-statistical approaches.

809 Note that, due to the different expressiveness of the languages underpinning
810 ILP and SRL methods w.r.t. those for the SW representations, the application

811 of such solutions is not straightforward. This problem has been considered since
812 the early works that apply machine learning methods to DL knowledge bases.
813 For instance, in [51], the authors have shown that there may be an exponential
814 blowup in knowledge base size and there may be some formulae without a coun-
815 terpart in DLs. Further issues have been discussed in [52] where the author
816 argues that *ad-hoc* solutions may avoid both exploring a larger search space
817 (represented by the set of all possible Horn clauses) and the limitations of the
818 complex reasoning services required by logic programming.

819 In order to better represent the inherent uncertainty related to the specific
820 semantics of the SW knowledge bases, the DST [21] offers an interesting al-
821 ternative, which explicitly considers the ignorance deriving from the inherent
822 incompleteness of the KBs and the availability of further evidence. Additionally,
823 the DST has been successfully integrated in various machine learning algorithms
824 to enhance the predictiveness of the models. For instance, DST primitives have
825 been integrated in the k -nearest neighbor algorithm [53], where each example
826 in the neighborhood is considered as a distinct source of evidence in favor of a
827 class that is subsequently combined through Dempster’s rule [21]. In the SW
828 context, a DL-compliant version of this approach has been proposed for solving
829 the class-membership prediction problem [54]. The DST has been integrated
830 also with algorithms for learning neural networks [55] and decision trees [56].
831 Indeed, the latter inspired our idea of evolving TDTs towards the ETDTs [7].
832 Differently from the original version of such a model (which is intended for
833 a propositional representation), the induction of an ETDT is guided by the
834 non-specificity measure whereas the original model considers also conflictual
835 evidences. In this paper we have extended our investigation considering further
836 total uncertainty measures.

837 The DST has been employed in the context of ensemble learning for pooling
838 the prediction coming from the weak learners [13, 15]. Various ensemble combi-
839 nation methods resort to *decision templates*, which are obtained by fitting, for
840 each classifier against each class, a mean vector (called *reference vector*). When
841 these models are employed, predictions are typically made by computing the

842 similarity between a decision profile of an unknown instance with the decision
843 templates. Unlike such approaches, the decision procedure employed with the
844 ETRFs combines the predictions returned in the form of BBAs. In this sense,
845 this procedure is similar to the one proposed in [57]: each classifier returns a
846 BBA that is combined by the meta-learner implementing a combination rule.
847 Again, ETRFs work on multi-relational representation language, similarly to
848 their original version, namely the *Terminological Random Forests* [10]. This
849 ensemble model, which represents a subtype of the First Order Logic Random
850 Forests [12] that is compliant with DLs, has been devised to tackle the problem
851 of class-imbalance in datasets drawn from Semantic Web knowledge bases (and
852 to overcome the limits of other solutions, such as those adopting sole sampling
853 methods [26]), which is an issue that had not been tackled before. A random
854 forest model for Semantic Web knowledge bases has been also proposed in [58]
855 but, unlike TRFs and ETRFs, the solution exploits only atomic concepts as
856 features.

857 One of the contributions of this paper concerns the adoption of a pruning
858 procedure for ETRFs, which mitigates some problems derived from the use of
859 many classifiers (e.g. the inefficiency in the prediction step) and can determine
860 a good forest size per learning problem. In general, the problem of determining
861 such number is still an open issue: even in the case of simpler representation
862 languages (attribute-value and propositional logic), there were only few works
863 that propose solutions which are often based on the use of statistical tests (e.g.
864 *McNeimar's test*) [59]. Instead, this number is a parameter whose value is
865 typically intended as user-provided [60]. Only in a recent work regarding the
866 application of random forests on data streams [61], the authors argued that the
867 ideal number of classifiers is strictly related to the number of class labels of the
868 dataset.

869 7. Conclusion and Extensions

870 We have proposed and extended a framework for inducing evidential ter-
871 minological decision trees and random forests, as developments of the termino-
872 logical decision trees and random forests, devised as solutions of the problem
873 of class-membership prediction for Semantic Web knowledge bases. Following
874 the lessons learned with previous versions, the new models tackle various short-
875 comings affecting the quality of the models, especially the cases of uncertain
876 classification and imbalanced datasets due to the inherent incompleteness of
877 the knowledge bases of interest. The resulting models combine predictions that
878 are represented as basic belief functions rather than votes, exploiting evidence
879 combination rules proposed in the context of the Dempster-Shafer Theory for
880 making the final decision. In addition, for evidential terminological random
881 forests, a strategy for optimizing the ensemble has been proposed.

882 Extensive experiments have been performed to assess the validity of the
883 proposed models, also considering datasets drawn from various Web ontolo-
884 gies, varying conditions and parameter settings, and in comparison with other
885 inductive models and learning strategies. The experiments have shown how
886 the proposed classification model can achieve a better predictiveness than the
887 previous versions of terminological decision trees and random forest. In various
888 cases, the results are better than the other learning systems. Moreover, the mod-
889 els tended to assign a definite membership yielding to induce a large number
890 of non logically derivable assertions whose correctness was assessed under the
891 Strong Disjointness Assumption [30].

892 Besides, the predictiveness of the evidential terminological decision trees
893 was found not to depend on the rule adopted for combining evidence while the
894 predictiveness of evidential terminological random forests was not affected by
895 the choice of either the forest size or the sampling rate. The standard deviation is
896 also lower than the one observed with the original TRFs. The evaluation showed
897 that the simplification procedure proposed to optimize the ensemble favors the
898 prediction of uncertain membership .

899 In the future, we plan to extend the method along various directions. One
900 regards considering an explicit *semi-supervised* learning approach for DL clas-
901 sifiers so to assign a definite membership to the uncertain examples. In this
902 case, it could be possible to devise solutions inspired from multi-view learning
903 approaches [62]. In addition, it can be interesting to investigate the effectiveness
904 of kernels derived from evidential random forests, as proposed in [63].

905 Further ensemble techniques and novel rules for combining the answers of the
906 weak learners could be employed. For example, weak learners can be induced
907 from subsets of training instances generated by means of a procedure based on
908 cross-validation rather than sampling with replacement. Further investigations
909 may concern the application of strategies aiming at the optimization of the
910 ensembles during the induction of the classifier rather than *ex post*, i.e. after the
911 training phase has been completed.

912 Finally, the (ensemble) methods could be naturally parallelized and the re-
913 sulting decision procedure based on induced models could be made available as
914 a service i.e. a non-standard inference service to complement standard query an-
915 swering or reasoning services. In this perspective, using specific frameworks such
916 as *Apache Spark*⁹ or GPUs may be an interesting alternative to be considered.

917 References

- 918 [1] A. Rettinger, U. Lössch, V. Tresp, C. d’Amato, N. Fanizzi, Mining the
919 Semantic Web - Statistical learning for next generation knowledge bases,
920 Data Min. Knowl. Discov. 24 (2012) 613–662.
- 921 [2] C. d’Amato, N. Fanizzi, B. Fazzinga, G. Gottlob, T. Lukasiewicz, Ontology-
922 based semantic search on the web and its combination with the power of
923 inductive reasoning, Ann. Math. Artif. Intell. 65 (2012) 83–121.
- 924 [3] N. Fanizzi, C. d’Amato, F. Esposito, Induction of concepts in web ontolo-
925 gies through terminological decision trees, in: J. Balcázar, et al. (Eds.),

⁹<http://spark.apache.org>

- 926 Proceedings of ECML/PKDD2010, volume 6321 of *LNAI*, Springer, 2010,
927 pp. 442–457.
- 928 [4] H. Blockeel, Top-down induction of first order logical decision trees, Ph.D.
929 thesis, Department of Computer Science, Katholieke Universiteit Leuven,
930 1998.
- 931 [5] H. Blockeel, L. De Raedt, Top-down induction of first-order logical decision
932 trees, *Artif. Intell.* 101 (1998) 285–297.
- 933 [6] F. Baader, D. Calvanese, D. McGuinness, D. Nardi, P. Patel-Schneider
934 (Eds.), *The Description Logic Handbook*, Cambridge University Press, 2nd
935 edition, 2007.
- 936 [7] G. Rizzo, C. d’Amato, N. Fanizzi, On the effectiveness of evidence-based
937 terminological decision trees, in: F. Esposito, et al. (Eds.), *Proceedings of*
938 *ISMIS 2015*, volume 9384, Springer, 2015, pp. 139–149.
- 939 [8] F. Smarandache, D. Han, A. Martin, Comparative study of contradiction
940 measures in the theory of belief functions, in: *Proceedings of FUSION*
941 2012, pp. 271–277.
- 942 [9] H. He, Y. Ma, *Imbalanced Learning: Foundations, Algorithms, and Appli-*
943 *cations*, Wiley-IEEE Press, 1st edition, 2013.
- 944 [10] G. Rizzo, C. d’Amato, N. Fanizzi, F. Esposito, Tackling the class-imbalance
945 learning problem in semantic web knowledge bases, in: K. Janowicz, et al.
946 (Eds.), *Proceedings of EKAW 2014*, volume 8876 of *LNCS*, Springer, 2014,
947 pp. 453–468.
- 948 [11] L. Breiman, Random forests, *Machine Learning* 45 (2001) 5–32.
- 949 [12] A. V. Assche, C. Vens, H. Blockeel, S. Dzeroski, First order random forests:
950 Learning relational classifiers with complex aggregates, *Machine Learning*
951 64 (2006) 149–182.

- 952 [13] L. Kuncheva, A theoretical study on six classifier fusion strategies, *Pattern*
953 *Analysis and Machine Intelligence*, IEEE Transactions on 24 (2002) 281–
954 286.
- 955 [14] L. Xu, A. Krzyzak, C. Suen, Methods of combining multiple classifiers and
956 their applications to handwriting recognition, *Systems, Man and Cyber-*
957 *netics*, IEEE Transactions on 22 (1992) 418–435.
- 958 [15] G. Rogova, Combining the results of several neural network classifiers, in:
959 R. Yager, L. Liu (Eds.), *Classic Works of the Dempster-Shafer Theory of*
960 *Belief Functions*, volume 219 of *Studies in Fuzziness and Soft Computing*,
961 Springer, 2008, pp. 683–692.
- 962 [16] X.-C. Yin, C. Yang, H.-W. Hao, Learning to diversify via weighted kernels
963 for classifier ensemble, *CoRR* abs/1406.1167 (2014).
- 964 [17] Y. Bi, J. Guan, D. Bell, The combination of multiple classifiers using an
965 evidential reasoning approach, *Artificial Intelligence* 172 (2008) 1731 –
966 1751.
- 967 [18] G. Rizzo, C. d’Amato, N. Fanizzi, F. Esposito, Inductive classification
968 through evidence-based models and their ensembles, in: F. Gandon, et al.
969 (Eds.), *Proceedings of ESWC 2015*, volume 9088 of *LNCS*, Springer, 2015,
970 pp. 418–433.
- 971 [19] K. Sentz, S. Ferson, Combination of evidence in Dempster-Shafer theory,
972 Technical Report, SANDIA, SAND2002-0835, 2002.
- 973 [20] J. Klir, *Uncertainty and Information*, Wiley, 2006.
- 974 [21] G. Shafer, *A Mathematical Theory of Evidence*, Princeton University Press,
975 Princeton, 1976.
- 976 [22] D. Dubois, H. Prade, On the combination of evidence in various mathemat-
977 ical frameworks, in: J. Flamm, T. Luisi (Eds.), *Reliability Data Collection*
978 *and Analysis*, volume 3 of *Eurocourses*, Springer, 1992, pp. 213–241.

- 979 [23] D. Dubois, H. Prade, A note on measures of specificity for fuzzy sets,
980 International Journal of General Systems 10 (1985) 279–283.
- 981 [24] U. Höhle, A general theory of fuzzy plausibility measures, Journal of
982 Mathematical Analysis and Applications 127 (1987) 346 – 364.
- 983 [25] R. R. Yager, Entropy and specificity in a mathematical theory of evidence,
984 in: R. Yager, L. Liu (Eds.), Classic Works of the Dempster-Shafer Theory
985 of Belief Functions, volume 219 of *Studies in Fuzziness and Soft Computing*,
986 Springer, 2008, pp. 291–310.
- 987 [26] H. He, E. A. Garcia, Learning from imbalanced data, IEEE Trans. on
988 Knowl. and Data Eng. 21 (2009) 1263–1284.
- 989 [27] J. Lehmann, S. Auer, L. Bühmann, S. Tramp, Class expression learning
990 for ontology engineering., J. Web Sem. (2011) 71–81.
- 991 [28] J. Shawe-Taylor, N. Cristianini, Kernel Methods for Pattern Analysis,
992 Cambridge University Press, New York, NY, USA, 2004.
- 993 [29] N. Fanizzi, C. d’Amato, F. Esposito, Induction of robust classifiers for web
994 ontologies through kernel machines, J. Web Sem. 11 (2012) 1–13.
- 995 [30] S. Schlobach, Debugging and semantic clarification by pinpointing, in:
996 Gómez-Pérez, et al. (Eds.), Proceedings of ESWC 2005, volume 3532 of
997 *LNCS*, Springer, 2005, pp. 226–240.
- 998 [31] K. Laskey, P. Costa, M. Kokar, T. Martin, T. Lukasiewicz, Uncertainty
999 Reasoning for the World Wide Web, Technical Report, 2008.
- 1000 [32] N. Fanizzi, C. d’Amato, F. Esposito, ReduCE: A reduced coulomb energy
1001 network method for approximate classification, in: L. Aroyo, et al. (Eds.),
1002 Proceedings of ESWC 2009, volume 5554 of *LNCS*, Springer, 2009, pp.
1003 323–337.
- 1004 [33] T. Hastie, R. Tibshirani, J. Friedman, The Elements of Statistical Learning,
1005 Springer, 2 edition, 2009.

- 1006 [34] N. Fanizzi, C. d'Amato, F. Esposito, DL-FOIL Concept Learning in De-
1007 scription Logics, in: Proceedings of ILP 2008, LNAI, Springer, 2008, pp.
1008 107–121.
- 1009 [35] V. Tripathy, A comparative study of multi-relational decision tree learning
1010 algorithm, International Journal of Scientific & Technology Research 2
1011 (2013).
- 1012 [36] H. A. Leiva, A multi-relational decision tree learning algorithm, Master's
1013 thesis, Department of Computer Science. Iowa State University, 2002.
- 1014 [37] D. Jeon, W. Kim, Development of semantic decision tree, in: The 3rd
1015 International Conference on Data Mining and Intelligent Information Tech-
1016 nology Applications, pp. 28–34.
- 1017 [38] L. Getoor, B. Taskar, Introduction to Statistical Relational Learning
1018 (Adaptive Computation and Machine Learning), The MIT Press, 2007.
- 1019 [39] L. De Raedt, P. Frasconi, K. Kersting, S. Muggleton (Eds.), Probabilistic
1020 Inductive Logic Programming: Theory and Applications, Springer, 2008.
- 1021 [40] K. Kersting, L. De Raedt, Basic principles of learning bayesian logic
1022 programs, in: L. De Raedt, et al. (Eds.), Probabilistic Inductive Logic
1023 Programming: Theory and Applications, volume 4911 of *LNCS*, Springer,
1024 2008, pp. 189–221.
- 1025 [41] M. Richardson, P. Domingos, Markov logic networks, *Mach. Learn.* 62
1026 (2006) 107–136.
- 1027 [42] P. M. Domingos, D. Lowd, S. Kok, H. Poon, M. Richardson, P. Singla, Just
1028 add weights: Markov logic for the semantic web, in: P. C. G. da Costa,
1029 et al. (Eds.), Uncertainty Reasoning for the Semantic Web I, ISWC Interna-
1030 tional Workshops, URSW 2005-2007, Revised Selected and Invited Papers,
1031 volume 5327 of *Lecture Notes in Computer Science*, Springer, 2008, pp.
1032 1–25.

- 1033 [43] P. Domingos, S. Kok, D. Lowd, H. Poon, M. Richardson, P. Singla, Markov
1034 Logic, Springer Berlin Heidelberg, Berlin, Heidelberg, pp. 92–117.
- 1035 [44] T. Khot, S. Natarajan, K. Kersting, J. Shavlik, Gradient-based boosting
1036 for statistical relational learning: the markov logic network and missing
1037 data cases, *Machine Learning* 100 (2015) 75–100.
- 1038 [45] S. Sarkhel, D. Venugopal, T. A. Pham, P. Singla, V. Gogate, Scalable
1039 training of markov logic networks using approximate counting, in: D. Schu-
1040 urmans, M. P. Wellman (Eds.), *Proceedings of the Thirtieth AAAI Con-
1041 ference on Artificial Intelligence, February 12-17, 2016, Phoenix, Arizona,
1042 USA.*, AAAI Press, 2016, pp. 1067–1073.
- 1043 [46] B. Ahmadi, K. Kersting, M. Mladenov, S. Natarajan, Exploiting symme-
1044 tries for scaling loopy belief propagation and relational training, *Machine
1045 Learning* 92 (2013) 91–132.
- 1046 [47] A. Kimmig, L. Mihalkova, L. Getoor, Lifted graphical models: A survey,
1047 *Mach. Learn.* 99 (2015) 1–45.
- 1048 [48] M. Nickel, V. Tresp, H. Kriegel, A three-way model for collective learning
1049 on multi-relational data, in: L. Getoor, T. Scheffer (Eds.), *Proceedings
1050 of the 28th International Conference on Machine Learning, ICML 2011,
1051 Bellevue, Washington, USA, June 28 - July 2, 2011, Omnipress, 2011, pp.
1052 809–816.*
- 1053 [49] M. Nickel, *Tensor factorization for relational learning*, Ph.D. thesis, Ludwig
1054 Maximilians University Munich, 2013.
- 1055 [50] A. Bordes, N. Usunier, A. García-Durán, J. Weston, O. Yakhnenko, Trans-
1056 lating embeddings for modeling multi-relational data, in: C. J. C. Burges,
1057 L. Bottou, Z. Ghahramani, K. Q. Weinberger (Eds.), *Advances in Neural
1058 Information Processing Systems 26: 27th Annual Conference on Neural
1059 Information Processing Systems 2013. Proceedings of a meeting held De-
1060 cember 5-8, 2013, Lake Tahoe, Nevada, United States.*, pp. 2787–2795.

- 1061 [51] L. Badea, S. H. Nienhuys-Cheng, A Refinement Operator for Description
1062 Logics, Springer Berlin Heidelberg, Berlin, Heidelberg, pp. 40–59.
- 1063 [52] J. Lehmann, Learning OWL Class Expressions, Ph.D. thesis, University of
1064 Leipzig, 2010. PhD in Computer Science, supervisors: Prof. Klaus-Peter
1065 Fährnich, Dr. Sören Auer.
- 1066 [53] T. Denoeux, A k-nearest neighbor classification rule based on dempster-
1067 shafer theory, *IEEE Transactions on Systems, Man, and Cybernetics* 25
1068 (1995) 804–813.
- 1069 [54] G. Rizzo, C. d’Amato, N. Fanizzi, F. Esposito, Assertion prediction with
1070 ontologies through evidence combination, in: F. Bobillo, et al. (Eds.),
1071 Uncertainty Reasoning for the Semantic Web II, volume 7123 of *LNAI*,
1072 Springer, 2013, pp. 282–299.
- 1073 [55] T. Denoeux, A neural network classifier based on dempster-shafer theory.,
1074 *IEEE Transactions on Systems, Man, and Cybernetics, Part A* 30 (2000)
1075 131–150.
- 1076 [56] N. Sutton-Charani, S. Destercke, T. Denoeux, Classification trees based on
1077 belief functions, in: T. Denoeux, M.-H. Masson (Eds.), *Belief Functions:*
1078 *Theory and Applications*, volume 164 of *Advances in Intelligent and Soft*
1079 *Computing*, Springer, 2012, pp. 77–84.
- 1080 [57] Y. Bi, J. Guan, D. Bell, The combination of multiple classifiers using an
1081 evidential reasoning approach., *Artif. Intell.* 172 (2008) 1731–1751.
- 1082 [58] D. Jeon, W. Kim, Random forest algorithm for linked data using a parallel
1083 processing environment, *IEICE Transactions on Information and Systems*
1084 E98.D (2015) 372–380.
- 1085 [59] P. Latinne, O. Debeir, C. Decaestecker, Limiting the number of trees in
1086 random forests, in: J. Kittler, F. Roli (Eds.), *Multiple Classifier Systems*,
1087 Second International Workshop, MCS 2001 Cambridge, UK, July 2-4, 2001,

- 1088 Proceedings, volume 2096 of *Lecture Notes in Computer Science*, Springer,
1089 2001, pp. 178–187.
- 1090 [60] T. M. Oshiro, P. S. Perez, J. A. Baranauskas, How many trees in a random
1091 forest?, in: Proceedings of the 8th International Conference on Machine
1092 Learning and Data Mining in Pattern Recognition, MLDM'12, Springer-
1093 Verlag, Berlin, Heidelberg, 2012, pp. 154–168.
- 1094 [61] H. R. Bonab, F. Can, A theoretical framework on the ideal number of
1095 classifiers for online ensembles in data streams, in: S. Mukhopadhyay,
1096 et al. (Eds.), Proceedings of the 25th ACM International Conference on
1097 Information and Knowledge Management, CIKM 2016, Indianapolis, IN,
1098 USA, October 24-28, 2016, ACM, 2016, pp. 2053–2056.
- 1099 [62] S. Sun, A survey of multi-view machine learning, *Neural Computing and*
1100 *Applications* 23 (2013) 2031–2038.
- 1101 [63] A. Davies, Z. Ghahramani, The random forest kernel and other kernels for
1102 big data from random partitions., *CoRR* abs/1402.4293 (2014).