

Multi-view overlapping clustering for the identification of the subject matter of legal judgments

Graziella De Martino^a, Gianvito Pio^{a,b}, Michelangelo Ceci^{a,b,c}

^a*Dept. of Computer Science, University of Bari Aldo Moro, Bari (Italy)*

^b*Big Data Lab. - National Interuniversity Consortium for Informatics, Rome (Italy)*

^c*Dept. of Knowledge Technologies, Jožef Stefan Institute, Ljubljana (Slovenia)*

Abstract

The legal field is generally burdened by paper-heavy activities, and the management of massive amounts of legal judgments without the adoption of computational tools may compromise the effectiveness and efficiency of administration processes. In this paper, we propose MOSTA, a novel unsupervised method to support the automated identification of groups of legal judgments with similar characteristics, with the goal of reducing the manual effort necessary for the management of legal judgments.

Methodologically, MOSTA learns two different embedding models for legal judgments. The first aims to represent the semantics of the textual content, while the second aims to represent co-citations of legal acts, also considering the granularity of the citations. Such representations are then fused through a multi-view approach based on an autoencoder, and the obtained representation is finally exploited by a novel overlapping clustering algorithm. The latter is an additional strong point of MOSTA, since, contrary to existing approaches, does not rely on additional input parameters that inherently influence the degree of overlap of the resulting clusters.

Our experiments, performed on three textual datasets, including a real-world legal dataset provided by EUR-Lex, proved that the proposed representation of cited legal acts, the adopted multi-view fusion strategy, and the novel overlapping clustering algorithm implemented in MOSTA provide a positive contribution to the quality of the identified clusters. Finally, MOSTA demonstrated to be able to outperform by a great margin existing complete solutions based on fine-tuned BERT embedding models and existing overlapping clustering algorithms.

Keywords: Multi-view overlapping clustering, Embedding, Legal judgments

1. Introduction

The *law* can be considered an ensemble of governance rules aiming to guarantee that the rights of the members of a community are not abused by other members, corporations, or authorities. These rules inherently define a framework for good governance that binds the society to ensure the safety and the justice of day-to-day activities. However, the legal sector is burdened by paper-heavy activities, and the manual management of massive amounts of legal documents may compromise the effectiveness and efficiency of administration processes. In this context, computational approaches, possibly based on Artificial Intelligence (AI) techniques, can support the transformation of slow, paper-based, processes into smart and efficient workflows, through the automated integration and analysis of massive amounts of data.

In the literature, we can find several works that proposed the application of AI techniques to solve different tasks in the legal field. For example, in [3], the application of solutions based on information technology in the legal field was deeply investigated. The authors first discussed how the law appears as a body of rules that can be represented and understood through automated reasoning, emphasizing the challenges raised by the presence of ambiguities and *open texture*. The authors also suggested the adoption of ontologies to represent crucial legal relationships and to support machine learning algorithms. Finally, the authors proposed LUIMA, an architecture based on the UIMA framework that proved to be able to perform the conceptual markup of legal documents considering the semantics. In [28] the authors proposed a tool that notifies lawyers and consumers about potentially unfair clauses listed in terms of service of online platforms. Mandal et al. [31] proposed a measure to assess the similarity between textual legal court documents to improve the accuracy and the scalability of legal document retrieval systems. Another relevant example is the work presented in [32], where the authors designed an automated data collection framework that detects eviction judgments issued by Dutch courts. The authors performed two experiments, where the emphasis was on locating eviction-related judgments and the resolution of the cases in the judgments, respectively.

Following this line of research, in this paper, we propose MOSTA (Multi-view Overlapping cluSTering of legAl judgments), a novel AI method that can identify groups of legal judgments with similar traits, possibly corre-

36 sponding to *subject matter(s)*¹, thus reducing the necessary human effort for
37 navigating, organizing, and classifying large quantities of legal judgments.
38 Note that even if subject matters could be considered as labels/categories
39 in supervised machine learning tasks, in real-world scenarios, labeled legal
40 judgments are scarcely available. This is the main motivation for which we
41 designed MOSTA as a novel unsupervised clustering approach. Specifically,
42 MOSTA falls in the category of *overlapping* clustering methods, that is, it is
43 able to assign each document to more than one cluster. The adoption of an
44 overlapping clustering approach in this scenario is motivated by the fact that
45 legal documents tend to be related to multiple subject matters [12, 30, 40],
46 and restricting legal judgments to belonging to a single cluster would lead to
47 disregard relevant secondary topics.

48 Another peculiarity of legal documents is that their complex semantics is
49 not entirely described by their textual content, but also by cited legal acts,
50 such as regulations, directives, decisions, recommendations, and opinions.
51 Moreover, legal citations can also guarantee that the judgment conclusion
52 is not based solely on the magistrate’s choice, but takes into account the
53 information conveyed by entrenched precedents. This aspect is particularly
54 important in legal systems based on Common law, which apply the *stare*
55 *decisis* principle. In such systems, the similarity of the scenario with respect
56 to that of precedents is exploited to push the decision towards a similar out-
57 come [27]. Note that although precedents play a less decisive role in Civil
58 law systems, they are frequently used to back, support and defend specific
59 outcomes, but also to show how a similar legal problem was previously dealt
60 with. Therefore, even if the textual content can properly be represented
61 by resorting to existing embedding techniques (e.g., BERT [15]), possibly
62 focused on the legal domain (e.g., LEGAL-BERT [10]), ignoring the infor-
63 mation conveyed by cited legal acts would lead to disregard relevant aspects
64 for the identification of the subject matters.

65 Although existing general-purpose overlapping clustering approaches can
66 overcome the limitation of a single cluster assignment, they usually require
67 additional input parameters, mainly to define the desired degree of overlap
68 [47]. Moreover, existing methods for document clustering are not able to
69 specifically take advantage of the complimentary information represented by

¹A *subject matter* denotes the substance of the arguments, reasoning and informal fallacies presented for consideration during a judgment hearing.

70 the cited legal acts, together with the textual content. Consequently, they
71 cannot accurately grasp the similarity between legal judgments.

72 In this context, the method MOSTA proposed in this paper solves all the
73 above-mentioned limitations. In particular, MOSTA is based on a multi-view
74 approach that fuses content-based embeddings with citation-based embed-
75 dings by means of a stacked autoencoder [5]. For the former, we adopt a
76 word embedding method able to consider the semantics of the textual con-
77 tent, as well as the contextual information. For the latter, we represent the
78 granularity of each citation (e.g., a whole act, an article, a sub-article, etc.)
79 through a tree-based structure, and exploit an embedding strategy based on
80 the similarity among trees. Note that MOSTA can work with any kind of
81 citations, both linking to precedents (typical of Common Law systems) and
82 linking to regulations, directives, and decisions (typical of Civil Law sys-
83 tems). Finally, MOSTA exploits a novel overlapping clustering method that
84 does not require additional input parameters, and is able to automatically
85 estimate the proper degree of overlap from data.

86 The rest of the paper is structured as follows. In Section 2 we describe
87 existing works related to the present paper, while in Section 3 we describe
88 in detail the proposed method MOSTA. In Section 4 we describe our exper-
89 imental evaluation, showing and discussing the obtained results. Finally, in
90 Section 5 we draw some conclusions and outline possible future work.

91 **2. Related Work**

92 In the following subsections, we briefly discuss existing approaches related
93 to the present paper. Specifically, we discuss existing clustering methods
94 applied in the legal field, and works that proposed multi-view document
95 clustering approaches, even if not specifically tailored for the legal field.

96 *2.1. Clustering of legal documents*

97 Most of the activities in the legal field are based on the management and
98 analysis of large amounts of textual documents. During the last years, the
99 increased availability of legal databases outlined new opportunities for auto-
100 mated data-driven approaches. In particular, in the literature we can find
101 several methods for cluster analysis, whose primary objective is the reduc-
102 tion of the complexity of repetitive tasks, by facilitating the navigation and
103 the organization of large collections of legal documents. A relevant exam-
104 ple is [11], where the authors applied clustering techniques to automatically

105 group case law petitions submitted to electronic trial systems. The authors
106 adapted the hard clustering algorithm initially proposed in [1], and intro-
107 duced the paradigm of *bag of terms and law references*. This paradigm is
108 based on a domain thesaurus to identify legal terms, and on regular expres-
109 sions (RE) to extract law references. Although this approach, similarly to
110 MOSTA, somehow considers citations, it treats them as textual words in the
111 bag, without properly considering their granularity. Moreover, the adopted
112 clustering method does not allow each document to fall into multiple clusters.

113 Lu et al. [30] proposed an overlapping clustering algorithm based on a
114 built-in topic segmentation approach that leverages legal metadata about
115 several types of legal documents. In addition to showing the scalability of
116 the proposed solution, the authors emphasized the ability to move from tra-
117 ditional lexical approaches toward the exploitation of topics, citations, and
118 click-stream data from behavior databases. However, the textual content is
119 represented through the classical bag-of-words model, with TF-IDF weigh-
120 ing, and the similarity among documents in terms of citations is based on
121 the Jaccard measure, without taking into account their granularity.

122 Conrad et al. [12] performed a comparative study between hard and
123 overlapping clustering solutions on three different legal datasets, using the
124 CLUTO clustering toolkit [50]. The results showed the effectiveness of over-
125 lapping and hierarchical clustering, in terms of both internal and external
126 quality measures, as well as in terms of the usefulness of the extracted clusters
127 for human legal experts. Similarly, Sabo et al. [40] explored the application
128 of approaches based on hard clustering (K-means and Affinity Propagation),
129 overlapping hierarchical clustering, and soft clustering (Lingo) to sparse nu-
130 merical vectors (obtained using the Bag of Words model) related to cases
131 dealing with airline service failure claims. The results showed the superiority
132 of hierarchical clustering in terms of entropy, purity, and legal experts' feed-
133 back. It is noteworthy that, in this case, possible overlaps among clusters
134 can occur only at different hierarchical levels, i.e., clusters can overlap simply
135 because of inclusive parent relationships. On the other hand, the considered
136 soft clustering solution requires a user-defined threshold to decide whether a
137 legal judgment belongs to a given cluster or not.

138 Existing general-purpose overlapping clustering approaches (e.g., [20]),
139 even if not specifically tailored for the legal field, can provide alternative
140 solutions if applied to a proper representation of legal documents. However,
141 analogously to soft clustering approaches, they require additional input pa-
142 rameters, that explicitly or implicitly influence the final degree of overlap.

143 An exception is represented by [47], which also proposes some strategies to
144 estimate the value of such additional parameters from data. For this reason,
145 in Section 4, we will consider it as a competitor with respect to the novel
146 clustering method implemented in MOSTA.

147 2.2. Multi-view document clustering

148 The need to take into account multiple perspectives/views of a document
149 could straightforwardly be satisfied by concatenating the features associated
150 with each different view. However, approaches based on feature concatena-
151 tion usually cannot differentiate the contribution provided by each view, and
152 could easily over-estimate the weight of a given view simply because it is rep-
153 resented by a high number of features. Therefore, in the literature, several
154 multi-view document clustering approaches have been proposed, that aim to
155 overcome the limitations of methods based on simple feature concatenation.

156 A relevant example is the work by Gao et al. [18], that extends the
157 information bottleneck algorithm to cluster web documents represented by
158 multiple distinct feature sets. Their experiments on two real datasets demon-
159 strated the effectiveness of the proposed approach, specifically when the views
160 represent the textual content, anchor texts, and URLs.

161 Other approaches are based on ensemble strategies. In particular, Kim
162 et al. [24] adopted an incremental algorithm to cluster multi-lingual docu-
163 ments, where each view provides a representation of documents in a different
164 language. In the first stage, the authors apply the Probabilistic Latent Se-
165 mantic Analysis (PLSA) [21] independently on each view, constraining each
166 clustering model to identify the same number of groups (topics). Then, they
167 identify the final clustering model such that documents falling in the same
168 group share similar patterns in terms of the probabilities returned by PLSA.
169 Wahid et al. [46] exploited a multi-objective optimization technique based on
170 the Non-Dominated Sorting Genetic Algorithm-II (NSGA-II) [14], aiming to
171 identify a clustering solution, among those returned by multiple clustering
172 methods applied to all the available views, that simultaneously minimizes
173 the number of obtained clusters, the number of words that are not in com-
174 mon among documents in the same cluster, and the inter-cluster similarity.
175 Hussain et al. [23] aggregated (by average) a cluster-based similarity matrix,
176 a pairwise similarity matrix, and an affinity matrix, computed through dif-
177 ferent approaches on the different views. A further clustering step is then
178 applied on the combined similarity matrix to obtain the final result. Finally,

179 Zamora and Sublime [48] combined clustering results obtained from different
180 views using an information theory model based on Kolmogorov complexity.

181 It is noteworthy that ensemble-based approaches (that work on the out-
182 put spaces) may suffer from similar issues with respect to approaches based
183 on feature concatenation (that work on the input spaces). Indeed, while in
184 the latter case each feature has the same importance, leading to possible
185 biases towards high-dimensional views, in ensemble-based approaches, each
186 view has the same importance, independently of the actual contribution it
187 provides. On the contrary, the approach implemented in MOSTA combines
188 the contribution provided by the features describing each view, without in-
189 troducing specific biases (see [2] for an overview on the effect of different
190 kinds of biases on the learned models).

191 Some other attempts to overcome this issue have been made in more
192 recent works. For example, Zhan et al. [49] proposed the multi-view graph-
193 regularized concept factorization (MVCF) method, based on concept factor-
194 ization. In addition to exploiting multi-view features, similarly to the system
195 SAIRUS [36], MVCF achieves superior clustering performances with respect
196 to previously-proposed methods by reducing the dimensionality of data and
197 by learning different weights for each view. Similarly, Bai et al. [4] designed a
198 deep neural network that learns a semantic mapping from a high-dimensional
199 to a low-dimensional feature space. In particular, the authors exploited a
200 neighbor-based autoencoder model and a cross-view autoencoder model to
201 involve neighbor-wise (within the same view) and view-wise complementary
202 information in the clustering process.

203 Although the above-mentioned methods can be considered as multi-view
204 clustering approaches, since they properly weigh the contribution provided
205 by different views, they are neither able to identify overlapping clusters nor
206 to properly capture the different granularities of legal citations we can find
207 in legal documents. In this respect, to the best of our knowledge, MOSTA
208 can be considered the first method that adopts a multi-view learning ap-
209 proach able to properly model both the textual content and citations of legal
210 acts, also considering their granularity, and that exploits a novel overlapping
211 clustering approach to identify their subject matters, without the need of
212 specifying additional parameters that influence the degree of overlap.

213 3. The proposed method MOSTA

214 Before describing the steps performed by our method MOSTA, we briefly
215 introduce some useful notation (see Appendix A for a compact view of all
216 the used symbols) and formally define the solved task. Let:

- 217 • J be a set of legal judgments, that also cite legal acts;
- 218 • k be the desired number of clusters, possibly representing legal subject
219 matters.

220 The task solved by MOSTA consists in the identification of k , possibly over-
221 lapping, clusters of the legal judgments J , taking into account *i*) the seman-
222 tics of their textual content and *ii*) the legal acts they cite, at different levels
223 of granularity (e.g., a whole act, an article, a sub-article, etc.). As per the
224 definition of overlapping clustering, each legal judgment $J_i \in J$ can possibly
225 be assigned to multiple clusters, representing the fact that it may be related
226 to multiple subject matters.

227 Our method consists of four main phases, namely:

- 228 1. **Embedding of the textual content of legal judgments**, that con-
229 sists in *i*) learning an embedding model from J , capable to represent
230 the *semantics of the textual content* of the judgments into a numerical
231 feature space, and *ii*) adopting the learned model to represent each
232 judgment $J_i \in J$ in the learned feature space.
- 233 2. **Embedding of the citations of legal judgments**, that consists in
234 *i*) learning an embedding model from J , capable of representing the
235 *co-citation network* (also considering the granularity of the citations)
236 of legal judgments towards legal acts into a numerical feature space,
237 and *ii*) adopting the learned model to represent each judgment $J_i \in J$
238 in the learned feature space.
- 239 3. **Multi-view embeddings fusion**, that is the construction of a fused,
240 multi-view representation for each judgment $J_i \in J$ through a stacked
241 autoencoder that exploits both the content-based and the citation-
242 based embeddings identified in phases 1 and 2.
- 243 4. **Identification of overlapping clusters of legal judgments**, that
244 consists in the adoption of a novel overlapping clustering approach,
245 that discovers k homogeneous groups of legal judgments according to
246 their fused embeddings, without requiring additional input parameters
247 to determine the degree of overlap.

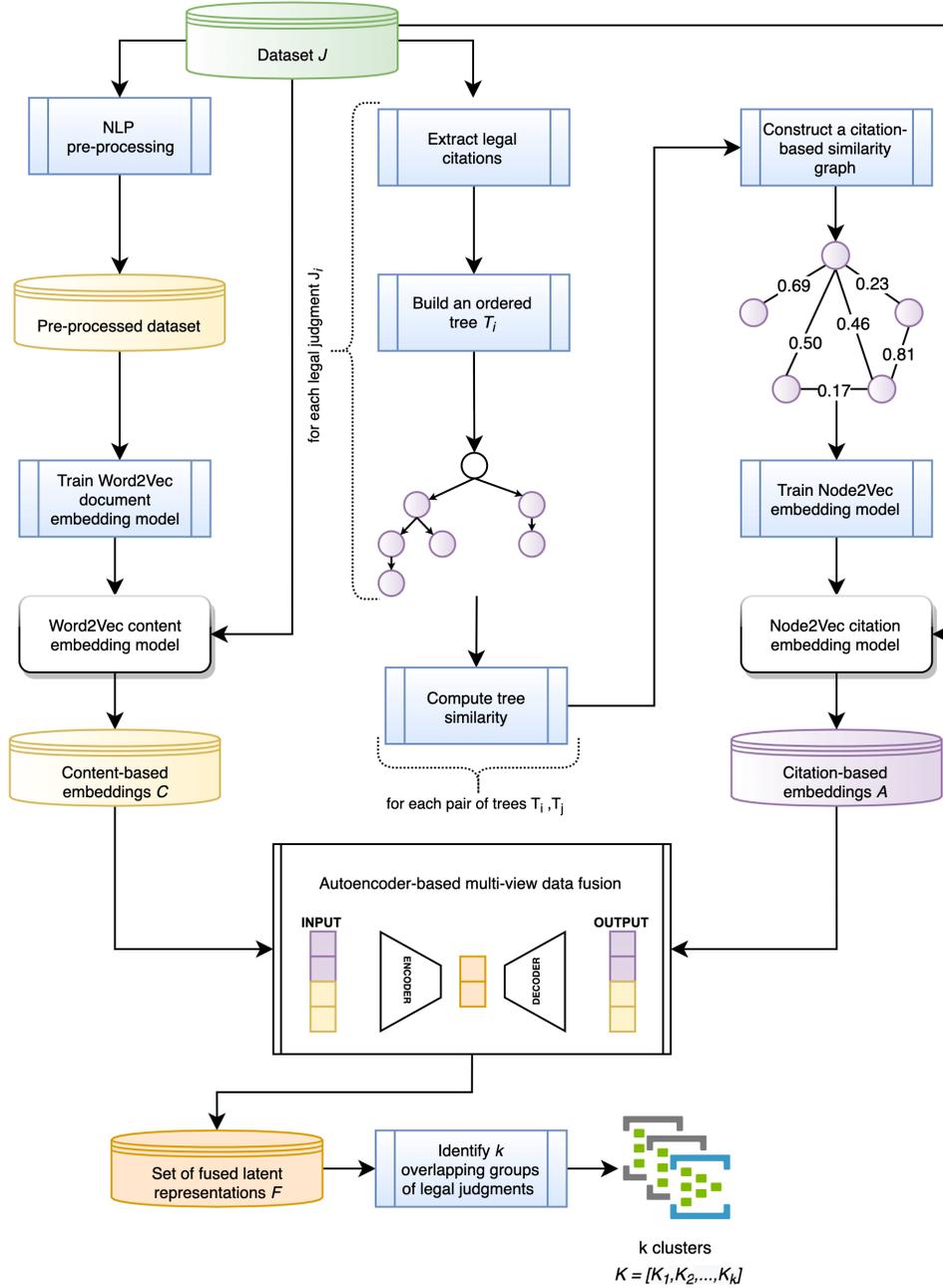


Figure 1: General workflow of the method MOSTA.

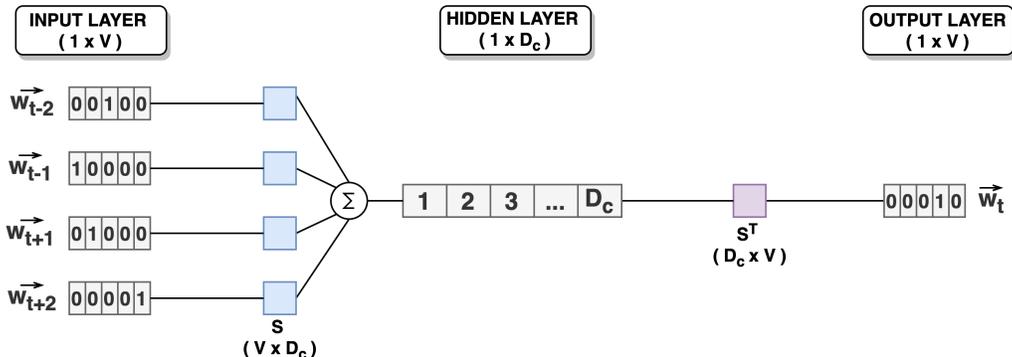


Figure 2: Graphical representation of the Word2Vec CBOW architecture. Note that there is only one matrix S , that is repeated multiple times only for explanatory purposes.

248 In the remainder of this section, we describe the approach followed by
 249 MOSTA to perform each phase, which is also globally depicted in Fig. 1.

250 3.1. Embedding of the textual content of legal judgments

251 In this section, we describe the steps followed by MOSTA for the represen-
 252 tation of the textual content of legal judgments in a numerical feature space.
 253 Initially, MOSTA adopts standard Natural Language Processing (NLP) [8]
 254 pre-processing techniques, namely, lowercasing, punctuation and digits re-
 255 moval, lemmatization, and removal of stopwords and rare words. Subse-
 256 quently, the pre-processed legal judgments are used to train an embedding
 257 model. In particular, MOSTA adopts the neural network (NN) architecture
 258 implemented in Word2Vec [33], given its proven superiority over traditional
 259 counting-based and other document-based embedding approaches, even in
 260 presence of noise in the data [31, 13, 26]. Word2Vec relies on two different
 261 shallow NN architectures, namely the Continuous-Bag-of-Words (CBOW)
 262 architecture and the Skip-gram (SG) architecture. Although both architec-
 263 tures are able to capture complex syntactic and semantic relationships among
 264 words, they adopt distinct learning processes. Specifically, CBOW aims to
 265 predict a target word from a surrounding context, while SG aims to predict
 266 the surrounding words of a given target word. The CBOW architecture is
 267 able to represent rare words more accurately, although it usually requires a
 268 slightly higher execution time than SG [37, 43]. Therefore, in MOSTA, we
 269 adopt the CBOW architecture, whose description is reported as follows.

270 Given a sequence of words $\langle w_{t-h}, \dots, w_t, \dots, w_{t+h} \rangle$ describing a target word

271 w_t and its context of size $2h$, the CBOW architecture takes as input the one-
 272 hot vector representation \vec{w}_i of size V for each context word w_i , where V is the
 273 size of the vocabulary observed in the set of legal judgments J . The learning
 274 phase aims to identify the optimal values for the matrix $S \in \mathbb{R}^{V \times D_C}$, where
 275 D_C represents the desired embedding dimensionality. The one-hot vector
 276 representation of each w_i is multiplied by S to obtain $2h$ vectors in \mathbb{R}^{D_C} .
 277 The hidden layer represents the embedding of the target word w_t obtained
 278 by aggregating the $2h$ vectors associated with the context words as follows:

$$\sum_{w_i \in \{w_{t-h}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+h}\}} \vec{w}_i \cdot S \quad (1)$$

279 The output layer is obtained by multiplying the embedding of the target
 280 word w_t by S^\top , and corresponds to the one-hot vector representation \vec{w}_t of
 281 the target word w_t . This means that the values of the matrix S are optimized
 282 so that the one-hot vector representation of the target word w_t is accurately
 283 reconstructed, given the one-hot vectors of the context words as inputs. The
 284 learned matrix S can therefore be used to embed any word into a numerical
 285 feature space of size D_C , given its context words.

286 Word2Vec naturally provides an embedding for each word. Therefore, in
 287 order to identify an embedding for the document in J , as suggested in [33],
 288 we adopt a mean aggregation strategy. The output of this phase is the set
 289 of embedded documents C , according to their textual content.

290 3.2. Embedding of the citations of legal judgments

291 During the redaction of legal documents, legal experts usually cite perti-
 292 nent legal acts, such as statutes, regulations, decisions, or directives [41]. A
 293 legal citation provides a direct link to a recognized source that *i*) references
 294 a legal act and/or a legal act section through which some conclusions are
 295 inferred; *ii*) supports the impartiality of the judgment, providing possible
 296 links to similar contexts and precedents.

297 Given the importance of legal act citations, in MOSTA we define an
 298 approach that extracts a set of citation-based embeddings A from the legal
 299 judgments J . The goal is to identify a complimentary representation, with
 300 respect to that based on the textual content, that takes into account co-cited
 301 legal acts, as well as the granularity of the citations.

302 In detail, for each legal judgment $J_i \in J$, MOSTA represents cited legal
 303 acts as an ordered tree T_i (see Fig. 3). Note that cited legal acts may already

Finally, **Article 11(1) and (2) of Regulation No 1954/2003** provides that, on the basis of the information to be communicated to the Commission by the Member States, the Commission is to submit to the Council a proposal for a Regulation fixing the maximum annual fishing effort for each Member State and for each area and fishery defined in Articles 3 and 6 and that the Council, acting by qualified majority on the proposal from the Commission, is to decide on that effort. In implementation of that provision, the Council adopted **Regulation No 1415/2004**.

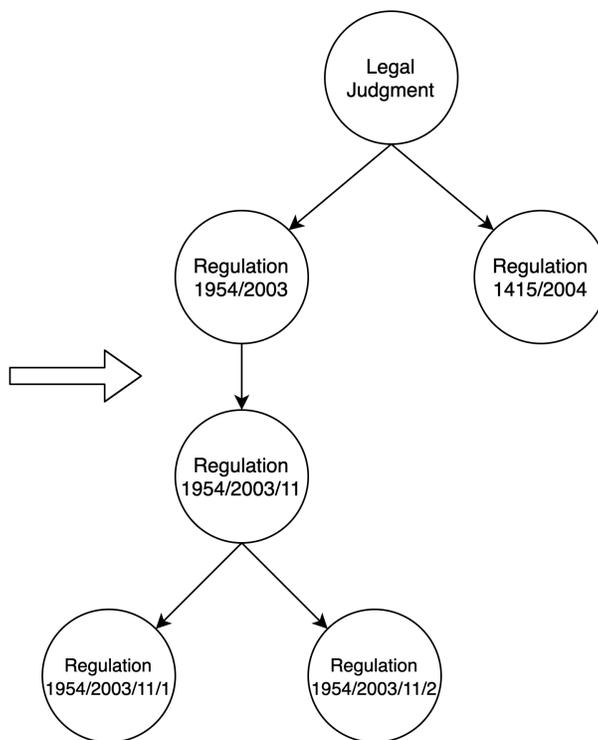


Figure 3: Representation of cited legal acts through an ordered tree.

304 be available as structured data in the dataset, or may need to be extracted,
 305 e.g., using Regular Expressions (RE). Since each legal system has its dis-
 306 tinctive characteristics and there is no uniformity across all jurisdictions, in
 307 Section 4.1, we define in detail the specific techniques used to extract legal
 308 citations from the dataset used in our experiments.

309 Once an ordered tree has been constructed for each $J_i \in J$, MOSTA
 310 computes the pairwise similarity between judgments. More formally, given
 311 two ordered trees T_i and T_j , extracted from the judgments $J_i \in J$ and $J_j \in J$,
 312 respectively, the tree similarity $s(T_i, T_j)$ is computed as:

$$s(T_i, T_j) = 1 - \frac{\delta(T_i, T_j)}{|T_i| + |T_j| - 2}, \quad (2)$$

313 where:

- 314 • $\delta(T_i, T_j)$ is the tree edit distance [35] defined as the minimum-cost se-
 315 quence of node edit operations, i.e., deletion, insertion, and relabeling

316 of nodes², needed to transform T_i into T_j ³;

- 317 • the factor $(|T_i| + |T_j| - 2)$, where $|\cdot|$ denotes the number of nodes of
318 a tree, corresponds to the maximum number of edit operations needed
319 to transform T_i into T_j , assuming that they are totally different trees.

320 To compute $\delta(T_i, T_j)$, MOSTA adopts the memory-efficient algorithm
321 APTED [35]. In Fig. 4, we report a step-by-step example of the compu-
322 tation of the tree edit distance, while in Fig. 5 we report multiple examples
323 of the similarity computed between different pairs of trees.

324 After computing the similarity between documents in terms of their cita-
325 tions, MOSTA builds a weighted graph $\mathcal{G} = (\mathcal{N}, \mathcal{E})$, where the set of nodes
326 \mathcal{N} corresponds to the judgments J , and each edge $\langle J_i, J_j \rangle \in \mathcal{E}$ represents the
327 fact that J_i and J_j co-cited some legal acts. Moreover, each edge $\langle J_i, J_j \rangle \in \mathcal{E}$
328 is associated with a weight corresponding to the similarity of their citations,
329 namely to $s(T_i, T_j)$, computed through Eq. (2).

330 Starting from such a weighted graph, we learn a numerical representation
331 for each node of the graph (i.e., for each judgment), where the new numerical
332 feature space aims to preserve the closeness relationships in the graph, also
333 according to the defined edge weights. In this way, the learned representa-
334 tion for a given judgment encodes the information about the fact that other
335 judgments co-cite the same legal acts, taking into account the granularity of
336 such co-citations thanks to the similarity measure defined in Eq. (2).

337 For the learning phase of the numerical representation from such a graph,
338 MOSTA exploits PecanPy [29], a memory-efficient implementation of the
339 method Node2Vec [19]. Node2Vec is a neural network architecture that learns
340 continuous feature representations for each node in a graph, by sampling
341 some representative nodes (in its neighborhood) following r 2^{nd} -order random
342 walks of fixed length l , biased by a hybrid Depth-First (DFS) / Breadth-First
343 (BFS) search approach. In particular, assuming that a given random walk
344 traverses the edge $\langle J_i, J_j \rangle$, the transition probability from J_j to the node
345 representing another judgment J_k , via the edge $\langle J_j, J_k \rangle$, is computed as

$$s(T_j, T_k) \cdot \beta(J_i, J_k) \tag{3}$$

²The cost of the relabeling operation is considered the double of the cost required for insertion or deletion operations, since it corresponds to a deletion of a node and to an insertion of a new node with a different label.

³Note that the considered node distance measure is symmetric. Therefore, the cost of transforming T_i into T_j is the same as that required to transform T_j into T_i .

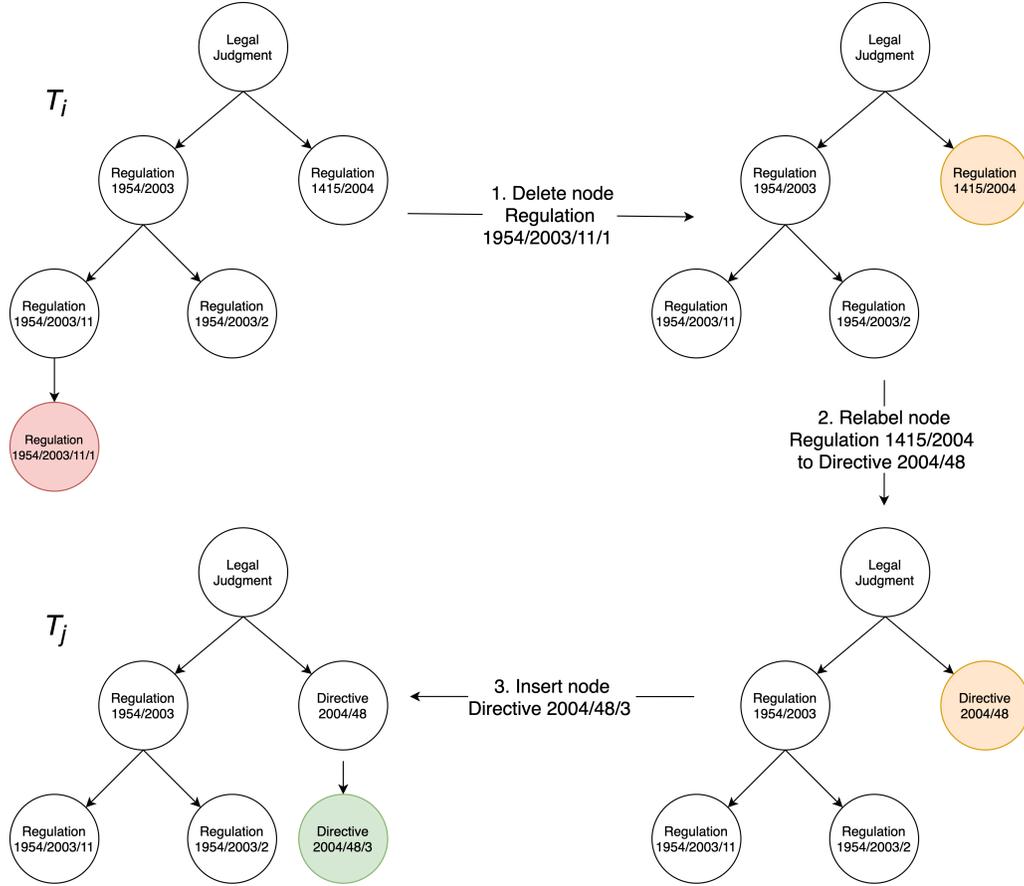


Figure 4: Graphical representation of the minimum-cost sequence of node edit operations needed to transform T_i into T_j . In the example, the distance between T_i and T_j is 4, which derives from the cost (1) for a node deletion operation (red) + the cost (2) for a node relabeling operation (orange) + the cost (1) for a node insertion operation (green).

346 where:

$$\beta(J_i, J_k) = \begin{cases} \frac{1}{p} & \text{if } g(J_i, J_k) = 0 \quad (\text{i.e., } J_i = J_k) \\ 1 & \text{if } g(J_i, J_k) = 1 \\ \frac{1}{q} & \text{if } g(J_i, J_k) = 2 \end{cases} \quad (4)$$

347 In Eq. (4), $g(J_i, J_k)$ is the distance (in terms of steps in the graph) between
 348 the nodes representing the judgments J_i and J_k ; p is a parameter that controls

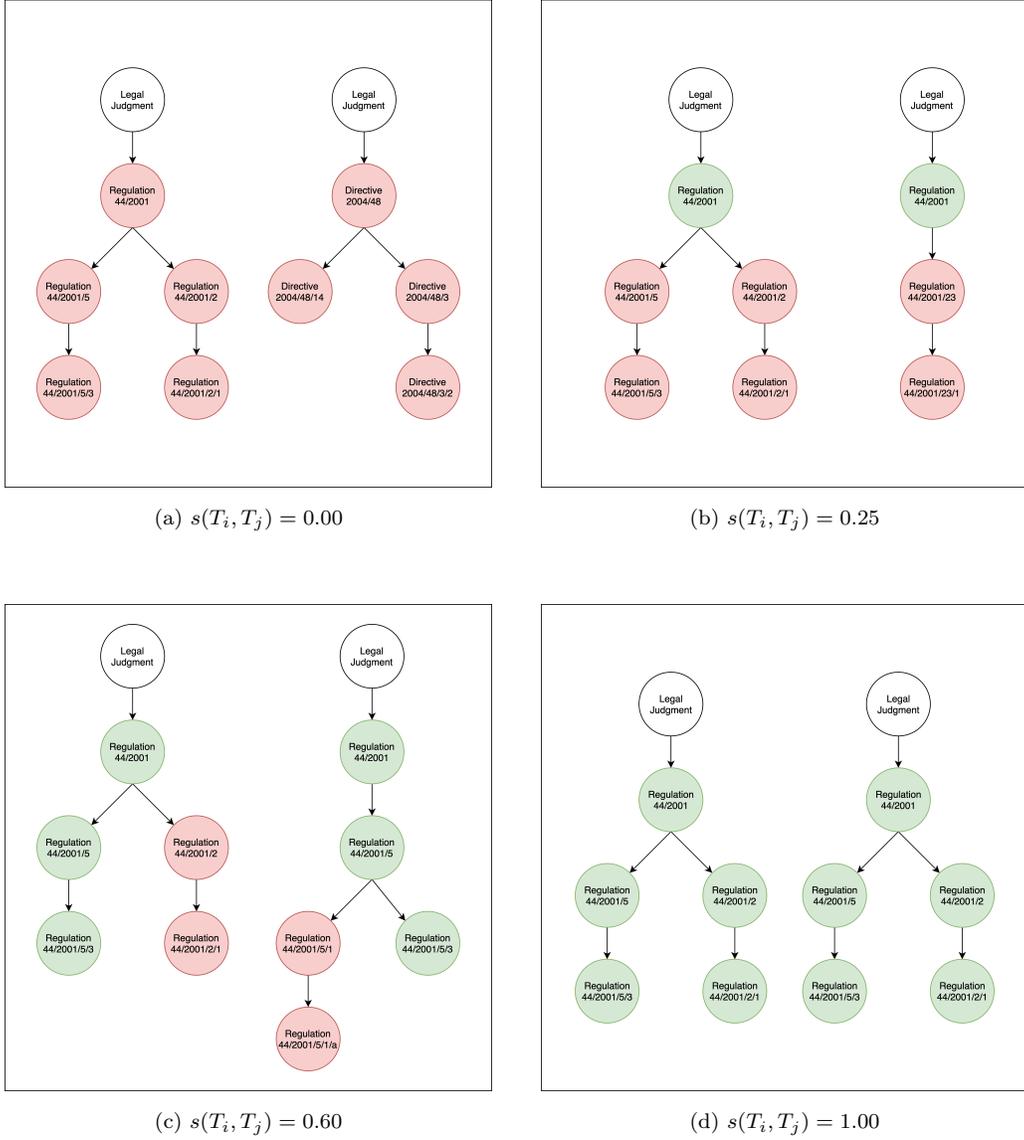


Figure 5: Examples of tree similarity scores computed between two ordered trees T_i, T_j . Green nodes represent matched citations, while red nodes represent differences.

349 the likelihood of immediately revisiting a node; q is a parameter that controls
350 how far the random walk should progress from J_i .

351 Subsequently, sampled random walks starting from each judgment are
352 considered as sequences of words representing its context, and are used to
353 learn a Word2Vec model. The embedding layer of this model is finally used
354 to extract the citation-based embeddings A for all the judgments J .

355 3.3. Multi-view embeddings fusion

356 The exploitation of multiple perspectives/views for the same units of
357 analysis has attracted increasing attention in the research community since,
358 when available, they can offer complimentary representations that may boost
359 the performance of the learned models. However, simple approaches, such
360 as feature concatenation, may introduce additional issues, namely feature
361 redundancy and collinearity [17], if the considered views are not completely
362 independent/orthogonal, and the curse of dimensionality, if the final number
363 of features is significantly higher than the available observations. As shown in
364 Section 2.2, more advanced approaches can be adopted, to properly capture
365 the contribution coming from the available views. In MOSTA, we adopt
366 an Autoencoder (AE) [5] to learn a low-dimensional fused representation
367 from the D_C -dimensional content-based embeddings C and from the D_A -
368 dimensional citation-based embeddings A .

369 An AE is an unsupervised feedforward neural network that learns a *com-*
370 *pressed* representation, such that the original data can be accurately recon-
371 structed. It comprises an *encoding* part, that maps the original input data
372 into the compressed space, and a *decoding* part, that reconstructs the original
373 data from its compressed version.

374 Methodologically, MOSTA initially concatenates content-based and
375 citation-based embeddings, leading to a feature vector in $\mathbb{R}^{D_C+D_A}$ for each
376 judgment. The input layer of the AE takes such a concatenated represen-
377 tation, which is compressed into a D_F -dimensional feature space, where
378 $D_F < D_C + D_A$. The specific architecture of the adopted AE is depicted
379 in Fig. 6. Note that, in general, multiple hidden layers can be defined in the
380 AE architecture before reaching the bottleneck layer that represents the final
381 embeddings. The choice of the number of additional hidden layers, as well
382 as of the number of their neurons, usually depends on the difference between
383 the input dimensionality ($D_C + D_A$, in our case) and the desired embedding
384 dimensionality (D_F , in our case).

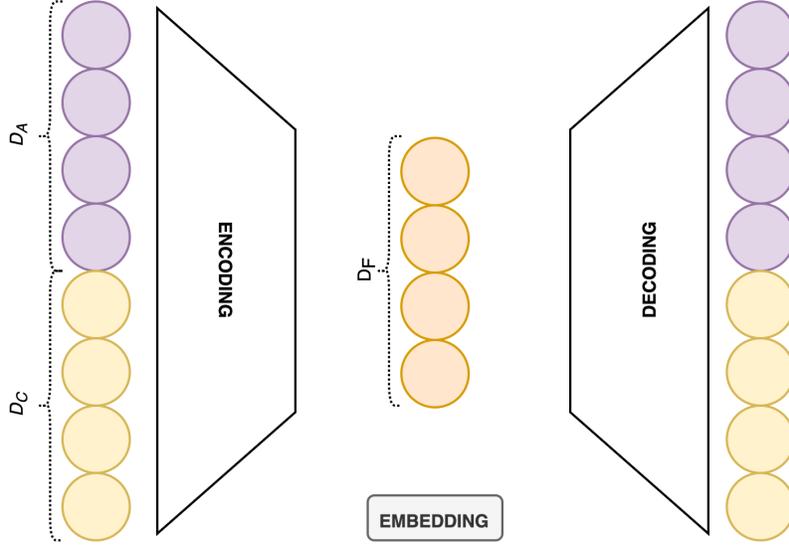


Figure 6: Graphical representation of the architecture of the adopted AE.

385 Note that the goal of the learning phase of the autoencoder is to opti-
 386 mize the weights of the neurons, such that the reconstruction errors, i.e.,
 387 the *loss* between the input and the output layer, is minimized. The loss is
 388 usually based on common measures like Root Mean Square Error (RMSE).
 389 In MOSTA, we adopt a different customized measure, that is able to provide
 390 different importance to the different sets of features belonging to each view.
 391 Specifically, we adopt a weighted variant of the RMSE, defined as follows:

$$\theta = \sqrt{\frac{1}{|J|} \sum_{J_i \in J} \gamma \cdot (\hat{x}_i - x_i)^2} \quad (5)$$

- 392 • x_i is the input $(D_C + D_A)$ -dimensional feature vector representing the
 393 judgment J_i ;
- 394 • \hat{x}_i is the $(D_C + D_A)$ -dimensional feature vector representing the judg-
 395 ment J_i , returned by the output layer of the AE;
- 396 • $\gamma = \lambda^{1 \times D_C} \oplus (1 - \lambda)^{1 \times D_A}$ defines the weights for the features coming
 397 from each view. If $\lambda = \frac{D_C}{D_C + D_A}$, θ corresponds to the standard RMSE.

398 Note that γ influences the importance given to each view in the computation
 399 of the loss function θ . Therefore, even if $\lambda = 0$ (resp., $\lambda = 1$) does not formally

400 mean that the AE discards the features of the content-based (resp., citation-
 401 based) embeddings, it implies that they are ignored in the computation of
 402 the loss function. In particular, when the loss function is required to ignore
 403 the features related to a specific view, such features would actually provide
 404 a negligible contribution to the obtained fused embeddings, meaning that
 405 we can consider the configuration $\lambda = 0$ (resp., $\lambda = 1$) equivalent to the
 406 scenario in which only the citation-based embeddings (resp., the content-
 407 based embeddings) are considered.

408 The result of this phase is the set of embedded judgments F , represented
 409 in a new D_F -dimensional feature space that fuses the contribution of the
 410 initial embeddings learned in the previous phases. The embedded judgments
 411 will be the input of the final clustering phase, that is described in the following
 412 subsection.

413 *3.4. Identification of overlapping clusters of legal judgments*

414 This subsection describes the novel clustering method that we imple-
 415 mented in MOSTA to identify k overlapping groups of legal judgments from
 416 F . A common approach adopted by existing overlapping clustering meth-
 417 ods, such as Neo K-Means [47], consists in the application of hard clustering
 418 solutions and in the assignment of additional clusters to each object, accord-
 419 ing to some criteria. However, as mentioned in Section 2.1, such criteria are
 420 usually based on a user-defined parameter that defines the degree of overlap,
 421 or the number of additional cluster assignments to perform. In MOSTA,
 422 we overcome this issue by adopting an approach based on outlier detection.
 423 Specifically, after applying a hard clustering method (i.e., k -means), MOSTA
 424 computes the Euclidean distance between each judgment and the centroid
 425 of each identified cluster. Assuming a Normal distribution of the distances,
 426 MOSTA identifies the judgment-cluster pairs, not already identified by the
 427 initial run of k -means, whose distance can be considered as an outlier. Specif-
 428 ically, following the $3\text{-}\sigma$ rule, MOSTA assigns a judgment to a given cluster
 429 if their distance is less than $d_{max} = \bar{d} - 3\sigma$ (see Fig. 7), where \bar{d} and σ are
 430 the average distance and the standard deviation of distances, respectively,
 431 between a judgment and a cluster centroid identified by k -means.

432 In Alg. 1, we report a pseudocode description of the clustering algorithm
 433 implemented in MOSTA. The algorithm starts by adopting the k -means clus-
 434 tering algorithm to partition F into k non-overlapping clusters (Alg. 1, line
 435 2). Then, the Euclidean distance between each judgment and each centroid of
 436 each cluster is computed (Alg. 1, line 3), in order to compute the mean \bar{d} and

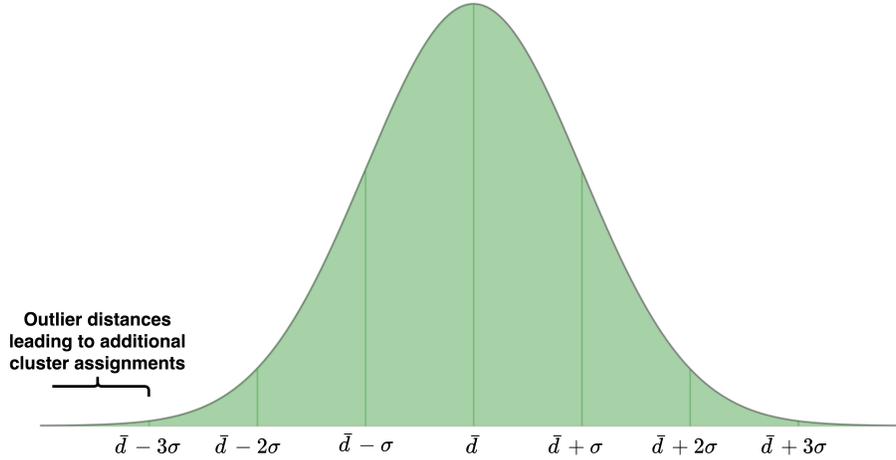


Figure 7: Representation of the outlier distances leading to additional cluster assignments.

437 the standard deviation σ (Alg. 1, line 4), and the threshold d_{max} to consider
 438 a distance between a judgement and a cluster as an outlier (Alg. 1, line 5).
 439 Finally, MOSTA performs additional cluster assignments when the observed
 440 judgment-cluster distance is less than the threshold d_{max} (Alg. 1, lines 6-9).
 441 We stress the fact that this strategy allows MOSTA to identify overlapping
 442 clusters by solely exploiting the observed distribution of distances, without
 443 imposing a pre-defined degree of overlap among clusters, or a pre-defined
 444 number of cluster assignments per judgement.

445 3.5. Time complexity analysis

446 In this subsection, we discuss the time complexity of the proposed method
 447 MOSTA, analyzing the time complexity of each phase, separately.

448 The first phase is the embedding of the textual content of legal documents,
 449 which corresponds to learning a Word2Vec model and to adopt it to embed
 450 all the documents. The time complexity of the training phase of Word2Vec
 451 is $O(|J| \cdot \log(V))$ [33], where V is the size of the vocabulary observed in the
 452 set of legal documents J . Once the learning phase of the Word2Vec model
 453 is completed, the embedding of each document requires $O(len \cdot 2h \cdot D_C \cdot V)$
 454 where len is the average number of words of a document, $2h$ is the size of
 455 the context and D_C is the dimensionality of the embedding. This complexity
 456 depends on a matrix multiplication between the input and the hidden layer
 457 of the Word2Vec architecture, performed for each word of the context of

458 each word of the document. Considering that h and D_C are constant values,
 459 the complexity of embedding all the documents can be approximated to
 460 $O(|J| \cdot len \cdot V)$.

461 The second phase of MOSTA requires the computation of a pairwise tree-
 462 based similarity between the documents, based on the citations. The com-
 463 putation of the similarity between two trees T_i and T_j (see Eq. (2)) has a
 464 time complexity $O(\max(|T_i|, |T_j|)^2)$ [35]. Considering that each tree is gen-
 465 erally very small, compared to the number of documents, the complexity of
 466 the pairwise tree-based similarity between the documents can be approxi-
 467 mated to $O(|J|^2)$. Moreover, the second phase requires learning an embed-
 468 ding model through Node2vec and its adoption to embed all the documents
 469 of the collection. Considering that Node2Vec is based on Word2Vec, where
 470 the vocabulary corresponds to the set of documents, its training complexity
 471 is $O(|J| \cdot \log|J|)$, while the cost for embedding all the documents is $O(|J|^2)$
 472 [33]. Therefore, the time complexity of the second phase can asymptotically
 473 be approximated to $O(|J|^2)$.

474 The third phase of MOSTA consists of learning a stacked autoencoder
 475 and using it to fuse the content-based and citation-based embeddings of each
 476 document. The time complexity of this step (see [6]) depends linearly on
 477 the number of training examples ($|J|$ in our case) and quadratically on the
 478 number of input features (in our case $(D_C + D_A)^2$). Therefore, considering
 479 that D_C and D_A are constant values and that $(D_C + D_A)^2 \ll |J|$, the time
 480 complexity of this phase can be approximated to $O(|J|)$.

481 Finally, we need to estimate the complexity of the clustering phase. The
 482 classical k -means algorithm linearly depends on both the number of examples
 483 and on the number of features. This means that, in our case, the time com-
 484 plexity of running k -means is $O(|J| \cdot D_F)$, where D_F is the fused embedding
 485 dimensionality. The identification of additional cluster assignments requires
 486 the computation of the distance between each document and each cluster
 487 centroid. Therefore, such a computation requires $O(|J| \cdot k \cdot D_F)$ operations.
 488 Finally, such distances are scanned once to identify those falling below the
 489 automatically generated threshold d_{max} . Since D_F and k are constants gen-
 490 erally much smaller than $|J|$, we can conclude that the time complexity of
 491 this phase can be approximated to $O(|J|)$.

492 Summing up the time complexity of each phase performed by MOSTA, we
 493 obtain $O(|J| \cdot len \cdot V) + O(|J|^2) + O(|J|) + O(|J|)$. If $(len \cdot V) \leq |J|$, the overall
 494 complexity of MOSTA can be approximated to $O(|J|^2)$, whereas if $(len \cdot V) >$
 495 $|J|$, the overall complexity of MOSTA can be approximated to $O(|J| \cdot len \cdot V)$.

496 Therefore, we can conclude that the time complexity of MOSTA is either
 497 quadratic in the number of documents to be processed or linear in the number
 498 of legal documents, in the average size of each document, and in the size of
 499 the vocabulary considered. Obviously, the worst-case analysis requires us to
 500 consider the highest time complexity among the two cases.

Algorithm 1: MOSTA overlapping clustering approach

Data:

- F : set of fused vector representations of the legal judgments J
- k : desired number of clusters to identify

Result:

- K : set of k overlapping clusters of legal judgments

```

1 begin
2   /* Identify  $k$  non-overlapping clusters through  $k$ -means */
3    $K \leftarrow k\text{-means}(F, k);$ 
4   /* Compute judgment-cluster pairwise distances */
5    $PD \leftarrow \text{computePairwiseDistances}(F, K);$ 
6   /* Compute mean and standard deviation of the distances */
7    $\bar{d} \leftarrow \frac{1}{|PD|} \cdot \sum_{d \in PD} d;$     $\sigma \leftarrow \sqrt{\frac{1}{|PD|} \sum_{d \in PD} (d - \bar{d})^2};$ 
8   /* Compute the threshold to consider a distance value as an outlier */
9    $d_{max} \leftarrow \bar{d} - 3 \cdot \sigma$ 
10  /* Identify overlapping clusters: perform additional judgment-cluster
11     assignments when their distance appears as an outlier */
12  foreach  $F_i \in F$  do
13    foreach  $K_j \in K$  do
14      if  $\text{distance}(F_i, K_j) < d_{max}$  and  $F_i \notin K_j$  then
15         $K_j \leftarrow K_j \cup \{F_i\};$ 
16      end
17    end
18  end
19  return  $K;$ 
20 end
```

502 4. Experiments

503 We performed the experiments along three different dimensions of analy-
504 sis. Specifically, we first evaluated the effectiveness of the proposed *overlap-*
505 *ping clustering method* implemented in MOSTA on three textual datasets,
506 in comparison with existing overlapping and soft clustering approaches.

507 Subsequently, we evaluated the effectiveness of the *multi-view fusion strat-*
508 *egy* adopted by MOSTA, and its ability to also capture the information con-
509 veyed by citations. This evaluation was performed on the EUR-Lex dataset⁴,
510 whose documents fall in the legal domain and provide both textual content
511 and citations.

512 Finally, on the same dataset, we compared the overall performance ex-
513 hibited by MOSTA with those achievable by complete competitor solutions
514 based on fine-tuned BERT embedding models and on the best overlapping
515 clustering method identified in the first phase of our experiments.

516 In the following subsections, we first detail the adopted datasets, the
517 competitor systems, the experimental setting and the evaluation measure.
518 Then, we show and discuss the obtained results for all the experiments.

519 4.1. The considered datasets

520 **EUR-Lex.** The first dataset that we considered in our experiments was
521 provided by EUR-Lex⁴. This dataset contains 4176 non-empty official public
522 EU legal judgments that were finalized between 2008 and 2018, categorized
523 in one or more *subject matters*⁵, that fall within the case-law sector and the
524 Court of Justice. In the dataset, we can find 133 distinct subject matters.

525 In order to build the set of citation-based embeddings A , we adopted
526 a custom strategy to extract citations from the dataset, since they were
527 not available as structured data. In particular, we reached EUR-Lex to
528 identify common rules adopted for citations in the legal judgments of this
529 specific dataset. Following their indications, we pre-processed the set of
530 judgments J by: *i*) lowercasing the text, *ii*) removing punctuation except for
531 the forward slash and the parenthesis (commonly used in citations), and *iii*)
532 removing stop words except for the word *of* (commonly used in citations).
533 Subsequently, we designed custom regular expressions (see Appendix B) to

⁴<https://eur-lex.europa.eu/homepage.html>

⁵For evaluation purposes, we discarded legal judgments not associated with any *subject matter* in the original dataset.

ID	Act Name	Article Number	Sub-Level 1	Sub-Level 2	Sub-Level 3	Sub-Level 4
62015CJ0005	Dir. 87/344	Dir. 87/344/4	Dir. 87/344/4/1	-	-	-
62015CJ0005	Dir. 87/344	-	-	-	-	-
62015CJ0005	Dir. 87/344	Dir. 87/344/4	Dir. 87/344/4/1	Dir. 87/344/4/1/a	-	-
62015CJ0005	Dir. 87/344	Dir. 87/344/3	Dir. 87/344/3/2	Dir. 87/344/3/2/c	-	-
62015CJ0005	Dir. 87/344	Dir. 87/344/3	Dir. 87/344/3/2	Dir. 87/344/3/2/a	-	-
...
62007CJ0416	Dir. 91/628	Dir. 91/628/5	Dir. 91/628-5/a	Dir. 91/628/5/a/1	Dir. 91/628/5/a/1/a	-
62007CJ0416	Reg. 806/2003	Reg. 806/2003/5	Reg. 806/2003/5/a	Reg. 806/2003/5/a/2	Reg. 806/2003/5/a/2/d	Reg. 806/2003/5/a/2/d/i
62007CJ0416	Dir. 90/425	-	-	-	-	-
62017CJ0530	Dec. 2015/143	Dec. 2015/143/2	Dec. 2015/143/2/1	-	-	-

Table 1: Examples of the structure of citations extracted from the judgments in the EUR-Lex dataset. *Dir.*, *Reg.* and *Dec.* are abbreviations of *Directive*, *Regulation* and *Decision*.

534 extract citations towards *Directives*, *Decisions*, and *Regulations*, following
535 the numbering rules for articles and sub-levels.

536 In Tab. 1, we show some examples of the structure of the citations as con-
537 firmed by EUR-Lex. A total of 36,116 unique legal citations were extracted,
538 leading, on average, to 8.65 cited acts per legal judgment.

539 Considering that this dataset specifically falls in the legal domain, and
540 that has both the textual content and citations, it has been exploited for all
541 the performed experiments, namely, *i*) for the evaluation of the effectiveness
542 of the proposed overlapping clustering approach implemented in MOSTA, *ii*)
543 for the evaluation of the performance of its fusion strategy, and *iii*) for the
544 comparison with existing complete solutions for the final task of identifying
545 the subject matters of legal documents.

546 **Reuters-21578.** This dataset consists of the train split of the ModHayes
547 Reuters-21578 subset⁶, which contains 9873 textual documents associated
548 with one or more topics, collected from the Reuters financial newswire service
549 during the 1987. The number of distinct topics in this dataset is 118.

550 **ArXiv.** This dataset consists of the train split of the arXiv dataset⁷, which
551 contains 4998 arXiv abstracts of submitted papers, associated with one or

⁶<https://huggingface.co/datasets/reuters21578>

⁷https://huggingface.co/datasets/arxiv_dataset

552 more system tags. In this dataset, we can find a total of 166 distinct tags.

553 For both Reuters-21578 and ArXiv, we pre-processed the text by *i*) low-
554 ercasing the text, *ii*) removing punctuation, and *iii*) removing stop words.
555 Considering that these datasets only contain textual content (i.e., with no
556 citations) and that each document is possibly associated with multiple top-
557 ics/tags, they were considered appropriate only for the first part of the ex-
558 periments, i.e., for the evaluation of the proposed overlapping clustering ap-
559 proach.

560 4.2. Experimental setting and competitor systems

561 In MOSTA, the embedding dimensionality of both the Word2Vec model
562 for the content-based embedding and the Node2Vec model for the citation-
563 based embedding was set to 256, i.e., $D_C = D_A = 256$, which is a pretty
564 standard value adopted for these architectures [45, 38]. The remaining pa-
565 rameters for Node2Vec were left to their default value, i.e., $p = 1$, $q = 1$,
566 $l = 80$ (length of random walks), and $r = 10$ (number of random walks).

567 For the evaluation of the clustering performance, we considered two com-
568 petitor algorithms. The first is Neo K-Means [47], which identifies over-
569 lapping clusters on the basis of a user-defined input threshold α . In Neo
570 K-Means, this parameter represents the average number of additional cluster
571 assignments per document. For the estimation of the optimal value of α , we
572 adopted the automatic strategy proposed in [47]. Moreover, we also evalu-
573 ated the results obtained when the optimal value of α is known a-priori, by
574 relying on the true number of cluster assignments in the dataset. Of course,
575 the results obtained in such a configuration are over-optimistic, since such
576 information is usually unknown in real scenarios. The second considered
577 competitor algorithm is Fuzzy C-Means (FCM) [7]. FCM is a soft-clustering
578 approach that returns the degree according to which each document belongs
579 to each cluster. In order to determine the cluster assignments, it requires a
580 user-defined threshold on such degrees. Since there is no automatic strategy
581 to determine such a threshold, we normalized the membership degrees in
582 $[0; 1]$ and collected the results with different thresholds, i.e., 0.3, 0.5, and 0.7.

583 As regards the parameter k , common to MOSTA and its competitors, we
584 run the experiments with different values following the rule of thumb, namely,
585 $k \in \{\sqrt{|J|}/2, \sqrt{|J|}, 2\sqrt{|J|}, 4\sqrt{|J|}, 8\sqrt{|J|}, 16\sqrt{|J|}\}$. The results with addi-
586 tional low values of k (e.g., $\sqrt{|J|}/4$, $\sqrt{|J|}/8$, and $\sqrt{|J|}/16$) are not reported,
587 since the obtained results appeared to be consistently worse with respect to

588 adopting higher values, for all the considered systems and parameter con-
589 figurations. On the other hands, we do not report the results with values
590 for k higher than $16\sqrt{|J|}$, because from $32\sqrt{|J|}$ the performance of MOSTA
591 naturally started to decrease since k was quickly degenerating to $|J|$ (note
592 that, in the EUR-Lex dataset, $|J| = 4176$, and $32\sqrt{|J|} = 2068$).

593 To specifically evaluate the effectiveness of the proposed multi-view fu-
594 sion strategy, we compared the results obtained by the AE implemented in
595 MOSTA with those achieved by other fusion approaches. The AE imple-
596 mented in MOSTA was structured with a simple 3-layers architecture with
597 only one hidden layer, corresponding to the bottleneck layer, with a dimen-
598 sionality of 256, namely, $D_F = 256$, and *sigmoid* as activation function. We
599 also evaluated the influence on the final results of the weight λ of the custom
600 loss function θ defined in Eq. (5). In particular, we performed the exper-
601 iments with $\lambda \in \{0.0, 0.1, 0.2, \dots, 0.8, 0.9, 1.0\}$. As competitor approaches,
602 we considered the simple concatenation of content-based and citation-based
603 embeddings, i.e., $A \oplus C$, and a feature weighting approach applied on such a
604 concatenation, which weighs each feature according to desired importance to
605 apply to its source view (i.e., content or citation embeddings). For this com-
606 petitor approach, we considered the following configurations of the weights for
607 content and citation embeddings, respectively: $\langle 0.1, 0.9 \rangle$, $\langle 0.3, 0.7 \rangle$, $\langle 0.7, 0.3 \rangle$,
608 and $\langle 0.9, 0.1 \rangle$. Intuitively, $\langle 0.1, 0.9 \rangle$ gives more importance to the content
609 than to citations, while $\langle 0.9, 0.1 \rangle$ does the opposite. Note that we did not
610 consider the configuration $\langle 0.5, 0.5 \rangle$ because it corresponds to the simple
611 concatenation approach $A \oplus C$, since it provides the same weight to all the
612 features.

613 Finally, as mentioned at the beginning of Sec. 4, we compared the results
614 achieved by MOSTA with those achievable by complete competitor solu-
615 tions, on the EUR-Lex dataset. Specifically, for the construction of content-
616 based embeddings, we considered pre-trained BERT-based models fine-tuned
617 for the legal field [10], namely LEGAL-BERT BASE (768-dimensional em-
618 beddings), LEGAL-BERT EURLEX (768-dimensional embeddings), and
619 LEGAL-BERT SMALL (512-dimensional embeddings). Note that LEGAL-
620 BERT EURLEX is specifically fine-tuned on the dataset adopted in this
621 evaluation, which, in principle, could provide it some advantages. Since
622 BERT-based models support the embedding of documents with maximum
623 512 tokens [16], we adopted two strategies [44, 34]: TS_1 , that preserves the
624 first 512 tokens of each legal judgment, and TS_2 that preserves the first and

625 the last part of each document, cutting off the middle part. For the clustering
626 phase, we considered Neo K-Means, since it provided the best results among
627 the competitors in the first part of our experiments aimed at evaluating the
628 clustering performance.

629 4.3. Evaluation measure

630 Since the datasets contain the true topics/subject matters assigned to
631 each document, as evaluation measure, we collected the F1-score averaged
632 over the clusters, computed after the identification of the best cluster-topic
633 matching through the Hungarian algorithm [25]. Therefore, for each cluster:

- 634 • a True Positive (TP) is a document/judgment that is labeled with the
635 topic/subject matter matched with the cluster;
- 636 • a False Positive (FP) is a document/judgment falling in the cluster
637 which is not labeled with the topic/subject matter matched with the
638 cluster;
- 639 • a True Negative (TN) is a document/judgment that did not fall in the
640 cluster and is not labeled with the topic/subject matter matched with
641 the cluster;
- 642 • a False Negative (FN) is a document/judgment that did not fall in the
643 cluster, but is labeled with the topic/subject matter matched with the
644 cluster.

645 Note that this evaluation setting is coherent with that usually adopted for
646 multi-label classification tasks [42].

647 The adoption of the average F1-score, instead of other measures like the
648 accuracy, is motivated by its ability to evaluate the quality of the result with-
649 out being biased by data unbalancing. Indeed, in the considered datasets, we
650 can notice a strong unbalancing (see Fig. 8). For the clustering task at hand,
651 the presence of unbalanced data corresponds to the fact that a few domi-
652 nant topics/subject matters, whose documents may be widely and unevenly
653 dispersed in the feature space, may partially obscure other topics/subject
654 matters that are less prominent and equally dispersed [9], making their mod-
655 eling by clustering algorithms much more difficult. This observation further
656 motivates the adoption of the F1-score as evaluation measure.

657 Note that, thanks to the availability of the ground truth in the datasets,
658 also in that specifically related to the legal field (EUR-Lex), we had the

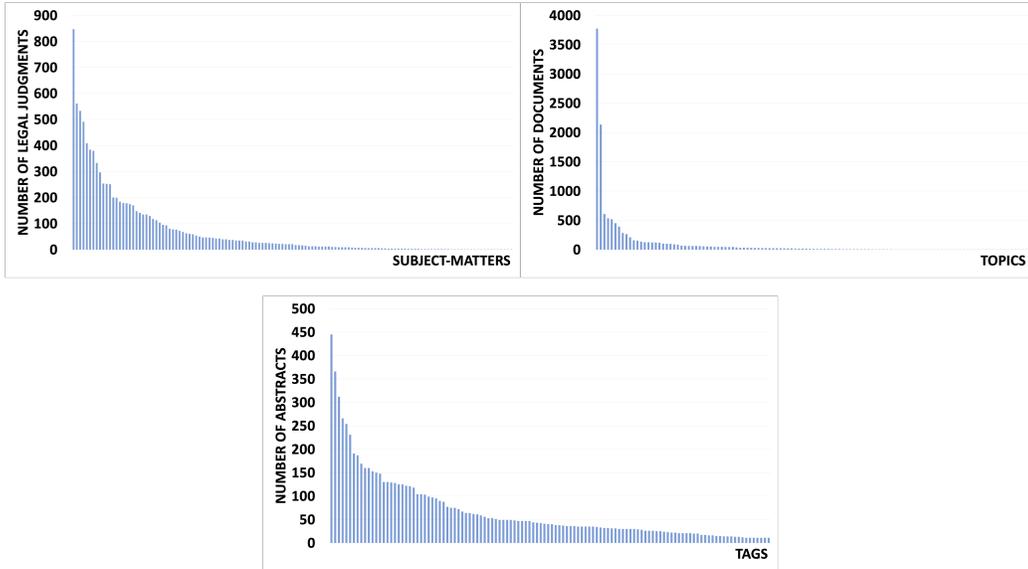


Figure 8: Number of legal judgments, documents and abstracts assigned to each subject matter, topic and tag, respectively, in the datasets EUR-Lex, Reuters and arXiv.

659 possibility to avoid the adoption of internal clustering quality measures, such
 660 as clustering agreement measures [39], since, when applied to overlapping
 661 clustering tasks, they tend to reward specific patterns in the resulting clusters
 662 (e.g., a low/high overlapping degree among clusters).

663 4.4. Results and Discussion

664 In Tab. 2, we report the F1-score results related to the evaluation of
 665 the overlapping clustering method implemented in MOSTA (see Sec. 3.4),
 666 applied only on the textual content of the three considered datasets.

667 We compared the results with those obtained by *i*) Neo K-Means with
 668 the automatic estimation of its parameter α , indicated as **N (est. α)**; *ii*)
 669 Neo K-Means with the optimal value of its parameter α , indicated as **N**
 670 **(opt. α)**; *iii*) Fuzzy C-Means, with different values of the threshold p applied
 671 on the membership degrees, indicated as **FCM_p**, with $p \in \{0.3, 0.5, 0.7\}$.
 672 The reported results refer to the F1-score obtained with different values of
 673 k . In the same table, we also report the average rank achieved by a given
 674 configuration, with respect to the clustering algorithm (last column of each
 675 sub-table) and k (last row of each sub-table).

676 Focusing on the value of k , we can observe higher F1-score results with

677 higher values of k . This is probably due to the high unbalancing in the
678 dataset (see Fig. 8), which makes the clustering algorithms more capable
679 of modeling the high amount of poorly represented subject matters in the
680 dataset when requiring a higher number of (thus, generally smaller) clusters.

681 Looking at the results obtained by different clustering algorithms, we
682 can easily conclude that MOSTA generally outperforms all the competitors.
683 This is also clear by observing the average ranks (last column of each sub-
684 table). The only case in which a competitor, i.e., Neo K-Means, is able to
685 compete with MOSTA is on the EUR-Lex dataset, but only when fed with
686 the ground value of its parameter α . The results obtained by FCM, with all
687 the considered values of its threshold, generally appear below those achieved
688 by MOSTA and Neo K-Means. Therefore, all the subsequent analyses have
689 been performed only considering these two algorithms.

690 In Tab. 3, we report the F1-score results obtained on the EUR-Lex
691 dataset, considering both the textual content and the citations, with different
692 values of λ for the multi-view fusion phase. In the same table, we also report
693 the average rank achieved by a given configuration, with respect to λ (last
694 column of each sub-table) and k (last row of each sub-table). As we can
695 observe from the table, also when using other algorithms for the clustering
696 phase, i.e., N (opt. α) and N (est. α), there is some influence coming from the
697 value of λ . Specifically, the best overall results were achieved with $\lambda = 0.1$
698 for MOSTA, $\lambda = 0.4$ for N (opt. α) and $\lambda = 0.2$ for N (est. α). This
699 result proves the usefulness of considering the information conveyed by the
700 citations in the multi-view fusion phase, irrespectively from the algorithm
701 adopted for the clustering phase. Therefore, citation-based embeddings can
702 be considered a useful complement to content-based embeddings, since they
703 positively contribute to the clustering results.

704 Overall, we can observe that the F1-score values obtained by MOSTA are
705 much higher than those obtained by N (est. α) and N (opt. α). In Tab. 4,
706 we make a direct comparison between MOSTA, N (est. α) and N (opt. α),
707 considering the best values of λ for each of them. As we can see from the
708 results, independently on the value of k , MOSTA consistently outperforms
709 N (est. α), and outperforms N (opt. α) in 4 out 6 cases, even if the latter
710 exploits the true value of α that, in principle, cannot be known a-priori. The
711 clear dominance of the clustering algorithm implemented in MOSTA (on
712 average, 10% higher F1-scores than N (opt. α) and 149% higher F1-scores
713 than N (est. α)), also confirmed by the average ranks (see the last row of
714 Tab. 4), confirms the effectiveness of the proposed outlier-based approach.

		k	$\sqrt{ J }/2$	$\sqrt{ J }$	$2\sqrt{ J }$	$4\sqrt{ J }$	$8\sqrt{ J }$	$16\sqrt{ J }$	AvgRank
		<i>alg.</i>							
EUR-LEX	MOSTA	<i>0.102</i>	<i>0.147</i>	<i>0.191</i>	0.225	0.254	0.278		1.50
	N (opt. α)	0.079	0.127	0.184	<i>0.242</i>	<i>0.283</i>	<i>0.314</i>		1.50
	N (est. α)	0.053	0.068	0.081	0.101	0.126	0.149		3.00
	$FCM_{0.3}$	0.029	0.034	0.037	0.039	0.043	0.046		4.42
	$FCM_{0.5}$	0.028	0.031	0.035	0.039	0.041	0.047		5.25
	$FCM_{0.7}$	0.028	0.034	0.035	0.037	0.041	0.046		5.33
AvgRank		6.00	5.00	4.00	3.00	2.00	1.00		
		k	$\sqrt{ J }/2$	$\sqrt{ J }$	$2\sqrt{ J }$	$4\sqrt{ J }$	$8\sqrt{ J }$	$16\sqrt{ J }$	AvgRank
		<i>alg.</i>							
REUTERS	MOSTA	<i>0.079</i>	<i>0.095</i>	<i>0.148</i>	<i>0.195</i>	<i>0.226</i>	<i>0.304</i>		1.00
	N (opt. α)	0.066	0.082	0.130	0.166	0.217	0.278		2.00
	N (est. α)	0.034	0.040	0.048	0.058	0.079	0.095		3.00
	$FCM_{0.3}$	0.024	0.025	0.025	0.025	0.025	0.025		6.00
	$FCM_{0.5}$	0.024	0.025	0.025	0.025	0.025	0.025		5.00
	$FCM_{0.7}$	0.024	0.025	0.025	0.025	0.025	0.025		4.00
AvgRank		6.00	5.00	3.33	2.83	2.17	1.67		
		k	$\sqrt{ J }/2$	$\sqrt{ J }$	$2\sqrt{ J }$	$4\sqrt{ J }$	$8\sqrt{ J }$	$16\sqrt{ J }$	AvgRank
		<i>alg.</i>							
ARXIV	MOSTA	<i>0.061</i>	<i>0.095</i>	<i>0.137</i>	<i>0.164</i>	<i>0.193</i>	0.209		1.17
	N (opt. α)	0.050	0.083	0.120	0.157	0.187	<i>0.224</i>		1.83
	N (est. α)	0.027	0.039	0.049	0.070	0.106	0.129		3.17
	$FCM_{0.3}$	0.018	0.025	0.030	0.033	0.033	0.033		5.17
	$FCM_{0.5}$	0.023	0.031	0.036	0.038	0.039	0.040		4.17
	$FCM_{0.7}$	0.028	0.015	0.020	0.018	0.016	0.017		5.50
AvgRank		5.17	5.17	3.67	2.83	2.33	1.83		

Table 2: F1-score results obtained on the textual content of the EUR-Lex, Reuters and arXiv datasets by the clustering algorithm implemented in MOSTA, Neo K-Means (opt. α), Neo K-Means (est. α) and Fuzzy C-Means with different thresholds applied on the membership degrees. Best column-wise results are emphasized with a gray background.

		k							
		λ	$\sqrt{ J }/2$	$\sqrt{ J }$	$2\sqrt{ J }$	$4\sqrt{ J }$	$8\sqrt{ J }$	$16\sqrt{ J }$	AvgRank
MOSTA		0.0	0.104	0.139	0.201	0.239	0.252	0.281	10.33
		0.1	0.124	<i>0.205</i>	0.250	0.280	<i>0.305</i>	0.301	2.00
		0.2	0.122	0.196	<i>0.261</i>	<i>0.286</i>	0.287	2.92	3.00
		0.3	0.121	0.194	0.254	0.263	0.298	0.294	3.50
		0.4	<i>0.127</i>	0.181	0.238	0.264	0.272	0.306	3.67
		0.5	0.116	0.181	0.228	0.256	0.277	<i>0.313</i>	4.17
		0.6	0.112	0.156	0.220	0.254	0.277	0.300	6.50
		0.7	0.120	0.164	0.216	0.253	0.270	0.296	6.50
		0.8	0.115	0.162	0.205	0.251	0.277	0.285	7.67
		0.9	0.115	0.161	0.216	0.246	0.272	0.293	8.00
		1.0	0.102	0.147	0.191	0.225	0.254	0.278	10.67
AvgRank			6.00	5.00	4.00	3.00	1.82	1.18	

		k							
		λ	$\sqrt{ J }/2$	$\sqrt{ J }$	$2\sqrt{ J }$	$4\sqrt{ J }$	$8\sqrt{ J }$	$16\sqrt{ J }$	AvgRank
N (opt. α)		0.0	0.057	0.097	0.158	0.174	0.188	0.245	11.00
		0.1	0.087	0.151	0.230	0.259	0.255	0.262	6.17
		0.2	0.085	0.159	0.211	0.239	0.286	0.282	6.83
		0.3	<i>0.098</i>	0.150	<i>0.238</i>	0.268	0.291	0.299	3.83
		0.4	0.096	<i>0.161</i>	0.238	<i>0.282</i>	0.290	0.323	2.50
		0.5	0.086	0.151	0.213	0.274	0.295	<i>0.324</i>	3.67
		0.6	0.092	0.142	0.206	0.266	<i>0.296</i>	<i>0.324</i>	4.17
		0.7	0.080	0.143	0.207	0.246	0.291	0.317	6.67
		0.8	0.091	0.151	0.194	0.258	0.294	0.305	5.50
		0.9	0.092	0.125	0.184	0.257	0.269	0.323	7.00
		1.0	0.079	0.127	0.184	0.242	0.283	0.314	8.67
AvgRank			6.00	5.00	4.00	2.91	2.00	1.09	

		k							
		λ	$\sqrt{ J }/2$	$\sqrt{ J }$	$2\sqrt{ J }$	$4\sqrt{ J }$	$8\sqrt{ J }$	$16\sqrt{ J }$	AvgRank
N (est. α)		0.0	0.031	0.043	0.055	0.078	0.115	0.137	8.83
		0.1	0.065	0.076	0.084	0.096	<i>0.144</i>	<i>0.149</i>	5.83
		0.2	0.066	0.081	0.090	0.097	0.136	0.142	5.00
		0.3	<i>0.070</i>	<i>0.085</i>	0.094	0.100	0.122	0.129	3.17
		0.4	0.067	0.084	0.095	<i>0.105</i>	0.111	0.114	3.17
		0.5	0.065	0.082	<i>0.096</i>	0.102	0.108	0.112	4.17
		0.6	0.064	0.081	0.093	0.103	0.106	0.109	6.17
		0.7	0.064	0.080	0.091	0.098	0.105	0.109	7.83
		0.8	0.064	0.079	0.091	0.099	0.105	0.110	7.33
		0.9	0.063	0.078	0.092	0.098	0.104	0.110	8.00
		1.0	0.053	0.068	0.081	0.101	0.126	<i>0.149</i>	6.50
AvgRank			6.00	5.00	3.91	3.09	2.00	1.00	

Table 3: F1-score results obtained on EUR-Lex (both content and citations) with different values of λ and k . The last column of each sub-table is the average rank of a given value of λ (by varying k), while the last row of each sub-table is the average rank of a given value of k (by varying λ) Best column-wise results are emphasized with a gray background.

	MOSTA	N	N
k		(opt. α)	(est. α)
$\sqrt{ J }/2$	<i>0.124</i>	0.096	0.067
$\sqrt{ J }$	<i>0.205</i>	0.161	0.084
$2\sqrt{ J }$	<i>0.250</i>	0.238	0.095
$4\sqrt{ J }$	0.280	<i>0.282</i>	0.105
$8\sqrt{ J }$	<i>0.305</i>	0.290	0.111
$16\sqrt{ J }$	0.301	<i>0.323</i>	0.114
AvgRank	1.33	1.67	3.00

Table 4: F1-score results obtained on EUR-Lex (both content and citations) by N (opt. α) and N (est. α), and by the clustering algorithm implemented in MOSTA, with the best fusion strategy (i.e., the AE implemented in MOSTA, as shown in Tab. 5), with their respective best value for λ , i.e., $\lambda = 0.1$ for MOSTA, $\lambda = 0.4$ for N (opt. α), and $\lambda = 0.4$ for N (est. α). Best row-wise results are emphasized with a gray background.

	MOSTA	C\oplusA	C\oplusA	C\oplusA	C\oplusA	C\oplusA
k			$\langle 0.1, 0.9 \rangle$	$\langle 0.3, 0.7 \rangle$	$\langle 0.7, 0.3 \rangle$	$\langle 0.9, 0.1 \rangle$
$\sqrt{ J }/2$	<i>0.124</i>	0.119	0.098	0.115	0.100	0.093
$\sqrt{ J }$	<i>0.205</i>	0.178	0.135	0.162	0.148	0.139
$2\sqrt{ J }$	<i>0.250</i>	0.214	0.200	0.214	0.194	0.185
$4\sqrt{ J }$	<i>0.280</i>	0.249	0.257	0.249	0.243	0.230
$8\sqrt{ J }$	<i>0.305</i>	0.262	0.276	0.269	0.270	0.260
$16\sqrt{ J }$	0.301	0.285	0.279	0.292	<i>0.309</i>	0.304
Avg. Rank	1.33	3.50	4.17	3.17	3.67	5.17

Table 5: F1-score results obtained on EUR-Lex (both content and citations) by the AE-based fusion strategy implemented in MOSTA ($\lambda = 0.1$) and by other fusion strategies based on the simple concatenation and on feature weighting, with different values of k . Best row-wise results are emphasized with a gray background.

715 In Tab. 5, we report the results of a further analysis aiming to specifi-
716 cally evaluate the contribution of the AE-based multi-view fusion strategy
717 implemented in MOSTA. In particular, we compare it with the concatena-
718 tion of the embeddings $C \oplus A$, as well as with an approach based on feature
719 weighting, considering different weights for each view (see Sec. 4 for details).

720 The results show that the proposed AE-based fusion strategy outperforms

721 the other considered techniques in almost all the situations (i.e., for almost
722 all the considered values of k). The influence of the weight on the feature
723 importance adopted for the considered competitor approach does not appear
724 to influence the results in a consistent way. In other words, determining the
725 best weight appears to be very challenging and dependent on the value of
726 k . For this specific analysis, we can conclude that the superiority of the AE-
727 based fusion strategy implemented in MOSTA is clear, and also confirmed by
728 the observed average ranks (see the last row of Tab. 5). These results confirm
729 that the proposed approach is able to significantly alleviate the issues possibly
730 introduced by the curse of dimensionality and to identify a fused feature
731 space that properly represents the complementary information conveyed by
732 the textual content and by cited legal acts.

733 Finally, in Tab. 6 we report the results of a comparison between the whole
734 method MOSTA and possible combinations of competitor systems that could
735 be adopted to solve the considered task on the EUR-Lex dataset. Specifi-
736 cally, as described in Sec. 4, we adopted different BERT-based embedding
737 models, and N (est. α) as the clustering algorithm. Note that, in this case,
738 a comparison with N (opt. α) would be totally unfair, since in real-world
739 scenarios, we cannot assume to know the true value of α . On the contrary,
740 both N (est. α) and MOSTA automatically identify the best estimate for
741 their parameters.

742 The F1-scores shown in Tab. 6 emphasize that MOSTA always outper-
743 forms all the competitors, independently on the adopted embedding model,
744 truncation strategy, and value of k . Indeed, MOSTA always ranks as the
745 first (best) method, in all the configurations (see the last row of Tab. 6).
746 On average, we can observe an improvement of 203%, 151% and 186% over
747 the results obtained when adopting LEGAL-BERT BASE, LEGAL-BERT
748 SMALL, and LEGAL-BERT EURLEX, respectively, as embedding models.
749 It is noteworthy that, among the competitor approaches adopted for the em-
750 bedding, LEGAL-BERT SMALL appears to be the best solution, even if not
751 specifically fine-tuned on the considered EUR-Lex dataset as LEGAL-BERT
752 EURLEX. This is probably due to the slightly lower number of features of
753 its embeddings (512 instead of 768), that alleviates the issues possibly intro-
754 duced by the curse of dimensionality. This observation further confirms the
755 appropriateness of the approach adopted by MOSTA.

756 Together with the specific analyses on the contribution provided by the
757 proposed overlapping clustering algorithm, by the citation-based embed-
758 dings, and by the multi-view AE-based fusion strategy, these final results

	LEGAL-BERT BASE		LEGAL-BERT SMALL		LEGAL-BERT EURLEX		MOSTA
k	TS ₁	TS ₂	TS ₁	TS ₂	TS ₁	TS ₂	
$\sqrt{ J }/2$	0.046	0.041	0.057	0.047	0.048	0.042	<i>0.124</i>
$\sqrt{ J }$	0.058	0.054	0.070	0.064	0.058	0.055	<i>0.205</i>
$2\sqrt{ J }$	0.074	0.064	0.080	0.083	0.071	0.069	<i>0.250</i>
$4\sqrt{ J }$	0.091	0.083	0.100	0.105	0.089	0.100	<i>0.280</i>
$8\sqrt{ J }$	0.109	0.109	0.133	0.148	0.127	0.122	<i>0.305</i>
$16\sqrt{ J }$	0.144	0.151	0.176	0.192	0.169	0.177	<i>0.301</i>
AvgRank	5.17	6.83	2.83	2.50	4.67	5.00	1.00

Table 6: F1-score results obtained on EUR-Lex (both content and citations) by MOSTA ($\lambda = 0.1$) and existing complete solutions, where the embedding is based on different BERT-based models, using the different truncation strategies TS_1 and TS_2 , and clustering is performed by N (est. α). Best row-wise results are emphasized with a gray background.

759 prove that the whole workflow implemented in MOSTA, that simultaneously
760 exploits the information conveyed by the textual content and by cited legal
761 acts, as well as its novel overlapping clustering method, can be considered
762 a precious tool for the unsupervised identification of the subject matters of
763 legal judgments.

764 5. Conclusions

765 In this paper, we proposed MOSTA, a novel method to identify groups
766 of legal judgments according to their characteristics. MOSTA is able to
767 identify a fused representation that considers both the textual content of
768 legal judgments and the legal acts they cite, properly taking into account the
769 granularity of the citations. Moreover, MOSTA adopts a novel overlapping
770 clustering method that does not require additional input parameters to define
771 the desired degree of cluster overlap, but automatically identifies additional
772 cluster assignments by exploiting an outlier-based strategy.

773 The specific evaluation of the performance of the proposed clustering
774 algorithm on three textual datasets proved that MOSTA is able to outper-
775 form Neo K-Means and Fuzzy C-Means, also considering different values of
776 their input parameters. Moreover, the experiments performed on a real le-
777 gal dataset provided by EUR-Lex emphasized that *i*) properly taking into
778 account citations can provide a positive contribution to the quality of the

779 identified clusters; *ii*) the proposed AE-based fusion strategy generally out-
780 performs concatenation-based approaches, including those that exploit fea-
781 ture weighting; *iii*) the clustering algorithm implemented in MOSTA out-
782 performs Neo K-Means, even when providing it with the optimal value of its
783 input parameter; *iv*) the whole method implemented in MOSTA outperforms
784 existing complete solutions based on the combination of pre-trained models
785 for document embedding and clustering.

786 For future work, we will take into account the aspects related to the ex-
787 plainability of the output, in order to make the clusters extracted by MOSTA
788 understandable and trustable [22]. Moreover, we will investigate the possibil-
789 ity to exploit the groups of legal judgments identified by MOSTA to provide
790 actual suggestions during the preparation of new legal judgments. In partic-
791 ular, we will explore the application of process mining techniques to clusters
792 of sequences of paragraphs to suggest the next paragraph to add to a legal
793 judgment under preparation.

Symbol	Description
J	A set of legal judgments
k	The number of clusters/groups of legal judgments to identify
<i>Embedding of the textual content of legal judgments</i>	
C	Content-based embeddings of the legal judgments J
D_C	Dimensionality of the content-based embeddings
w_i, \vec{w}_i	A context word and its one-hot vector representation
w_t, \vec{w}_t	A target word and its one-hot vector representation
h	Size of the context window
V	Size of the vocabulary observed in the set of legal judgments J
S	Weight matrix learned by the Word2Vec model
<i>Embedding of the citations of legal judgments</i>	
A	Citation-based embeddings of the legal judgments J
D_A	Dimensionality of the citation-based embeddings
T_i	An ordered tree representing the citations of the document J_i
$s(T_i, T_j)$	Tree similarity between the ordered trees T_i and T_j
$\delta(T_i, T_j)$	Tree edit distance between the ordered trees T_i and T_j
$\mathcal{G} = (N, E)$	A weighted graph. \mathcal{N} = legal judgments J ; \mathcal{E} = co-citations of legal acts
r, l	Number and length of Node2Vec random walks for each node
$\beta(J_i, J_k)$	The function defining the likelihood to reach the node J_k starting from J_i
$g(J_i, J_k)$	The distance (i.e., number of steps) between J_i and J_k in the graph
p, q	Node2Vec parameters to bias random walks
<i>Multi-view embeddings fusion</i>	
F	Fused, compressed, embeddings
D_F	The dimensionality of the fused latent representation
θ	The loss function adopted in the AE
λ	Importance of content-based embeddings in the AE loss
γ	The vector of weights used by the AE loss, based on the parameter λ
<i>Identification of overlapping clusters of legal judgments</i>	
k	Number of overlapping clusters of judgments to identify
K	Set of overlapping clusters of legal judgments
\bar{d}, σ	The mean and the standard deviation of the judgment-cluster distances

795 **Appendix B. Regular expressions for the extraction of citations**

796 In the following, we report the Regular Expressions adopted to extract the
797 citations from the legal judgments of the EUR-Lex dataset:

```
798 1: (?<!of\s)(council\s)*(?<!of\s)council\s)(regulation|decision|directive  
799 ↪ )((\s\((cfsp|ec|ecsc|eec|eu|euratom|jha|op_dat%pro)\))*\s\d+/\d  
800 ↪ +/(cfsp|ec|ecsc|eec|eu|euratom|jha|op_datpro))*+((\s(-| | | ))  
801 ↪ *(\s((articles?)|(arts?))(\s\d+)%)+((\s\d+)|([a-z])|(\(\w+\))|(\s  
802 ↪ s\(\w+\)))*)*,  
803  
804 2: (((articles?)|(arts?))\s\d+([a-z]*)+((\s\d+)|([a-z])|(\(\w+\))|(\s  
805 ↪ \(\w+\)))*)*\s(office)*\s(council\s)*(regulation|decision|directive)((\s  
806 ↪ s\((cfsp|ec|ecsc|eec|eu|euratom|jha|op_datpro)\))*\s\d+/\d+/(  
807 ↪ cfsp|ec|ecsc|eec|eu|euratom|jha|op_datpro))*))
```

808 **Availability:** The system, the dataset and all the results are available at:
809 https://osf.io/a9jm2/?view_only=471428680ce5483abc358fa17a99ad5f.

810 **Acknowledgements:** This work was partially supported by the project
811 FAIR - Future AI Research (PE00000013), Spoke 6 - Symbiotic AI
812 (CUP H97G22000210007), under the NRRP MUR program funded by
813 the NextGenerationEU. The authors also acknowledge the support of the
814 project “START UPP – Modelli, sistemi e competenze per l’implementazione
815 dell’ufficio per il processo” (CUP H29J22000390006) funded by the Italian
816 Minister of Justice.

817 **References**

- 818 [1] Aggarwal, C.C., Gates, S.C., Yu, P.S., 2004. On using partial supervision
819 for text categorization. *IEEE Transactions on Knowledge and Data*
820 *Engineering* 16, 245–255.
- 821 [2] Alessa, A., Jianlong, Z., Kevin, T., Fang, C., Andreas, H., 2022. Fairness
822 and Explanation in AI-Informed Decision Making. *Machine Learning*
823 *and Knowledge Extraction* 4, 556–579.
- 824 [3] Ashley, K.D., 2017. *Artificial Intelligence and Legal Analytics: New*
825 *Tools for Law Practice in the Digital Age*. Cambridge University Press.
- 826 [4] Bai, R., Huang, R., Chen, Y., Qin, Y., 2021. Deep multi-view document
827 clustering with enhanced semantic embedding. *Inf. Sci.* 564, 273–287.

- 828 [5] Ballard, D.H., 1987. Modular learning in neural networks, in: Proceed-
829 ings of the 6th National Conference on Artificial Intelligence. Seattle,
830 WA, USA, July 1987, Morgan Kaufmann. pp. 279–284.
- 831 [6] Barracchia, E.P., Pio, G., Bifet, A., Gomes, H.M., Pfahringer, B., Ceci,
832 M., 2022. LP-ROBIN: Link prediction in dynamic networks exploiting
833 incremental node embedding. *Information Sciences* 606, 702–721.
- 834 [7] Bezdek, J.C., Ehrlich, R., Full, W., 1984. Fcm: The fuzzy c-means
835 clustering algorithm. *Computers & Geosciences* 10, 191–203.
- 836 [8] Bird, S., Klein, E., Loper, E., 2009. *Natural Language Processing with
837 Python*. O’Reilly.
- 838 [9] Carrington, A.M., Fieguth, P.W., Qazi, H., Holzinger, A., Chen, H.H.,
839 Mayr, F., Manuel, D.G., 2020. A new concordant partial AUC and par-
840 tial c statistic for imbalanced data in the evaluation of machine learning
841 algorithms. *BMC Medical Informatics Decis. Mak.* 20, 4.
- 842 [10] Chalkidis, I., Fergadiotis, M., Malakasiotis, P., Aletras, N., Androu-
843 sopoulos, I., 2020. LEGAL-BERT: The muppets straight out of law
844 school, in: *Findings of the Association for Computational Linguistics:
845 EMNLP 2020*, Association for Computational Linguistics, Online. pp.
846 2898–2904.
- 847 [11] de Colla Furquim, L.O., de Lima, V.L.S., 2012. Clustering and cate-
848 gorization of brazilian portuguese legal documents, in: *PROPOR 2012*,
849 Coimbra, Portugal, April 17-20, 2012. *Proceedings*, pp. 272–283.
- 850 [12] Conrad, J.G., Al-Kofahi, K., Zhao, Y., Karypis, G., 2005. Effective doc-
851 ument clustering for large heterogeneous law firm collections, in: *Pro-
852 ceedings of the 10th International Conference on Artificial Intelligence
853 and Law*, p. 177–187.
- 854 [13] De Martino, G., Pio, G., Ceci, M., 2022. PRILJ: an efficient two-step
855 method based on embedding and clustering for the identification of reg-
856 ularities in legal case judgments. *Artificial Intelligence and Law* 30,
857 359–390.

- 858 [14] Deb, K., Agrawal, S., Pratap, A., Meyarivan, T., 2002. A fast and elitist
859 multiobjective genetic algorithm: NSGA-II. *IEEE Trans. Evol. Comput.*
860 6, 182–197.
- 861 [15] Devlin, J., Chang, M., Lee, K., Toutanova, K., 2018. BERT: pre-training
862 of deep bidirectional transformers for language understanding. *CoRR*
863 abs/1810.04805.
- 864 [16] Devlin, J., Chang, M.W., Lee, K., Toutanova, K., 2019. BERT: Pre-
865 training of deep bidirectional transformers for language understanding,
866 in: *ACL 2019, Association for Computational Linguistics*. pp. 4171–
867 4186.
- 868 [17] Draper, N., Smith, H., 1966. *Applied regression analysis*. Wiley series
869 in probability and mathematical statistics, Wiley.
- 870 [18] Gao, Y., Gu, S., Xia, L., Fei, Y., 2006. Web document clustering with
871 multi-view information bottleneck, in: *CIMCA 2006, International Con-
872 ference on Intelligent Agents, Web Technologies and Internet Commerce
873 (IAWTIC 2006)*, IEEE Computer Society. p. 148.
- 874 [19] Grover, A., Leskovec, J., 2016. node2vec: Scalable feature learning
875 for networks, in: *Proceedings of the 22nd ACM SIGKDD International
876 Conference on Knowledge Discovery and Data Mining, San Francisco,
877 CA, USA, August 13-17, 2016*, ACM. pp. 855–864.
- 878 [20] Hess, S., Pio, G., Hochstenbach, M.E., Ceci, M., 2021. BROCCOLI:
879 overlapping and outlier-robust biclustering through proximal stochastic
880 gradient descent. *Data Min. Knowl. Discov.* 35, 2542–2576.
- 881 [21] Hofmann, T., 2001. Unsupervised learning by probabilistic latent se-
882 mantic analysis. *Mach. Learn.* 42, 177–196.
- 883 [22] Holzinger, A., 2021. The next frontier: AI we can really trust, in:
884 *Machine Learning and Principles and Practice of Knowledge Discovery
885 in Databases - International Workshops of ECML PKDD 2021, Virtual
886 Event, September 13-17, 2021, Proceedings, Part I*, Springer. pp. 427–
887 440.
- 888 [23] Hussain, S.F., Mushtaq, M., Halim, Z., 2014. Multi-view document
889 clustering via ensemble method. *J. Intell. Inf. Syst.* 43, 81–99.

- 890 [24] Kim, Y., Amini, M., Goutte, C., Gallinari, P., 2010. Multi-view cluster-
891 ing of multilingual documents, in: ACM SIGIR 2010, Geneva, Switzer-
892 land, July 19-23, 2010, ACM. pp. 821–822.
- 893 [25] Kuhn, H.W., 2010. The hungarian method for the assignment problem,
894 in: 50 Years of Integer Programming 1958-2008 - From the Early Years
895 to the State-of-the-Art. Springer, pp. 29–47.
- 896 [26] Kumar, A., Makhija, P., Gupta, A., 2020. Noisy text data: Achilles’
897 heel of BERT, in: Proceedings of the Sixth Workshop on Noisy User-
898 generated Text, W-NUT@EMNLP 2020 Online, November 19, 2020,
899 Association for Computational Linguistics. pp. 16–21.
- 900 [27] Lewis, S., 2021. Precedent and the Rule of Law. *Oxford Journal of Legal*
901 *Studies* 41, 873–898.
- 902 [28] Lippi, M., Pałka, P., Contissa, G., Lagioia, F., Micklitz, H.W., Sartor,
903 G., Torroni, P., 2019. Claudette: an automated detector of potentially
904 unfair clauses in online terms of service. *Artificial Intelligence and Law*
905 27, 117–139.
- 906 [29] Liu, R., Krishnan, A., 2021. PecanPy: a fast, efficient and parallelized
907 Python implementation of node2vec. *Bioinformatics* 37, 3377–3379.
- 908 [30] Lu, Q., Conrad, J.G., Al-Kofahi, K., Keenan, W., 2011. Legal document
909 clustering with built-in topic segmentation, in: CIKM 2011, Glasgow,
910 United Kingdom, October 24-28, 2011, ACM. pp. 383–392.
- 911 [31] Mandal, A., Ghosh, K., Ghosh, S., Mandal, S., 2021. Unsupervised
912 approaches for measuring textual similarity between legal court case
913 reports. *Artificial Intelligence and Law* 29, 417–451.
- 914 [32] Medvedeva, M., Dam, T., Wieling, M., Vols, M., 2021. Automatically
915 identifying eviction cases and outcomes within case law of dutch courts
916 of first instance, in: *Legal Knowledge and Information Systems - JURIX*
917 *2021: The Thirty-fourth Annual Conference*, Vilnius, Lithuania, 8-10
918 December 2021, IOS Press. pp. 13–22.
- 919 [33] Mikolov, T., Sutskever, I., Chen, K., Corrado, G., Dean, J., 2013. Dis-
920 tributed representations of words and phrases and their composition-
921 ality, in: *NIPS 2013*, Curran Associates Inc., Red Hook, NY, USA. p.
922 3111–3119.

- 923 [34] Mutasodirin, M.A., Prasojo, R.E., 2021. Investigating Text Shortening
924 Strategy in BERT: Truncation vs Summarization, in: (ICAC SIS) 2021,
925 pp. 1–5.
- 926 [35] Pawlik, M., Augsten, N., 2016. Tree edit distance: Robust and memory-
927 efficient. *Inf. Syst.* 56, 157–173.
- 928 [36] Pellicani, A., Pio, G., Redavid, D., Ceci, M., 2023. SAIRUS: Spatially-
929 aware identification of risky users in social networks. *Information Fusion*
930 92, 435–449.
- 931 [37] Pennington, J., Socher, R., Manning, C.D., 2014. Glove: Global vectors
932 for word representation, in: *EMNLP 2014*, October 25-29, 2014, Doha,
933 Qatar, pp. 1532–1543.
- 934 [38] Qiao, Y., Zhang, B., Zhang, W., 2020. Malware classification method
935 based on word vector of bytes and multilayer perception, in: *2020 IEEE*
936 *International Conference on Communications, ICC 2020*, Dublin, Ire-
937 land, June 7-11, 2020, pp. 1–6.
- 938 [39] Rabbany, R., Zaïane, O.R., 2015. Generalization of clustering agree-
939 ments and distances for overlapping clusters and network communities.
940 *Data Mining and Knowledge Discovery* 29, 1458–1485.
- 941 [40] Sabo, I.C., Pont, T.R.D., Wilton, P.E.V., Rover, A.J., Hübner, J.F.,
942 2022. Clustering of brazilian legal judgments about failures in air trans-
943 port service: an evaluation of different approaches. *Artificial Intelligence*
944 *and Law* 30, 21–57.
- 945 [41] Sloan, A.E., 2018. *Basic legal research: Tools and strategies*. Wolters
946 Kluwer.
- 947 [42] Song, D., Vold, A., Madan, K., Schilder, F., 2022. Multi-label legal
948 document classification: A deep learning-based approach with label-
949 attention and domain-specific pre-training. *Information Systems* 106,
950 101718.
- 951 [43] Stratos, K., Collins, M., Hsu, D.J., 2015. Model-based word embeddings
952 from decompositions of count matrices, in: *ACL 2015*, July 26-31, 2015,
953 Beijing, China, Volume 1: Long Papers, pp. 1282–1291.

- 954 [44] Sun, C., Qiu, X., Xu, Y., Huang, X., 2019. How to fine-tune BERT
955 for text classification?, in: CCL 2019, Kunming, China, October 18-20,
956 2019, Proceedings, Springer. pp. 194–206.
- 957 [45] Tissier, J., Gravier, C., Habrard, A., 2019. Near-lossless binarization of
958 word embeddings, in: AAAI 2019, Honolulu, Hawaii, USA, January 27
959 - February 1, 2019, pp. 7104–7111.
- 960 [46] Wahid, A., Gao, X., Andreae, P., 2014. Multi-view clustering of web
961 documents using multi-objective genetic algorithm, in: CEC 2014, Bei-
962 jing, China, July 6-11, 2014, IEEE. pp. 2625–2632.
- 963 [47] Whang, J.J., Dhillon, I.S., Gleich, D.F., 2015. Non-exhaustive, overlap-
964 ping k -means, in: Proceedings of the 2015 SIAM International Confer-
965 ence on Data Mining, Vancouver, BC, Canada, April 30 - May 2, 2015,
966 SIAM. pp. 936–944.
- 967 [48] Zamora, J., Sublime, J., 2020. A new information theory based cluster-
968 ing fusion method for multi-view representations of text documents, in:
969 SCSM 2020, Held as Part of HCII 2020, Copenhagen, Denmark, July
970 19-24, 2020, Proceedings, Part I, Springer. pp. 156–167.
- 971 [49] Zhan, K., Shi, J., Wang, J., Tian, F., 2017. Graph-regularized con-
972 cept factorization for multi-view document clustering. *J. Vis. Commun.*
973 *Image Represent.* 48, 411–418.
- 974 [50] Zhao, Y., Karypis, G., Fayyad, U., 2005. Hierarchical clustering algo-
975 rithms for document datasets. *Data Min. Knowl. Discov.* 10, 141–168.