

Adebayo Kolawole John / Luigi Di Caro / Guido Boella

Annotating Legal Documents with Ontology Concepts

This paper describes a task of semantic labeling. The idea exploits ontology in providing a fine-grained conceptual document annotation. The proposed system performs conceptual tagging for efficient information filtering. The paper presents a promising solution. The proposed task has several applications such as granular information filtering of legal texts, text summarization and information extraction among others and has been evaluated on the task of conceptual tagging of semantic segments in text with promising result.

Category: Articles

Region: Italy

Field of law: Legal Informatics

Collection: Conference Proceedings IRIS 2016

Citation: Adebayo Kolawole John / Luigi Di Caro / Guido Boella, Annotating Legal Documents with Ontology Concepts, in: Jusletter IT IRIS

Contents

1. Introduction
2. Background and Related Work
3. Research Motivation
4. Methodology
 - 4.1. Concept Analysis
 - 4.2. Concept Zoning for Text Segmentation
 - 4.3. Concept-Document Mapping
5. Evaluation
6. Conclusion
7. References

1. Introduction

[Rz 1] The increasing use of computer, coupled with the growth of internet and emerging powerful database technologies implies that data accumulates in an unprecedented manner. Fortunately, huge amount of data now readily exist in electronic form otherwise called Electronically Stored Information (ESI). With the rising influence of Electronic Discovery (eDiscovery), the field of law (especially in terms of litigation) has tremendously benefitted from the availability and growth of ESI by allowing huge documents to be tendered in law courts for cross examination without the documents being in their physical form [1]. Just like the saying goes, «big data means big headache». The «big headache» in eDiscovery is technically a scaled-up Information Retrieval (IR) task in which manual classification of documents into either being relevant or producible or not for a civil, criminal or regulatory matter under litigation, is automated [2].

[Rz 2] In this paper, we pursue a form of fine-grained information retrieval on legal text, while we are not generally there yet; we have nevertheless developed a solution amenable to the realities of legal norm intricacies. The goal of this research is not to develop an eDiscovery system but rather a system whose output can further improve accuracy of eDiscovery systems as well as other tasks by giving more structure to legal texts as well as performing Semantic Annotation (SA) on legal texts for improved search facilities.

[Rz 3] We describe our idea in tagging legal documents, already divided into semantically coherent blocks with specific concept(s) (from a pool of concepts from Eurovoc¹) that describes the meaning of the content of the block. We opine that dividing documents into semantically coherent units and labeling each unit with its inferred concept can aid fine-grained information retrieval. To buttress this, we introduce the idea of semantic scopes to documents. This could be local or global scopes, describing how significant a concept is in representing the meaning of a document. The most significant concept being the global scope concept and the others are the local scope concepts for each document. Information search can thereafter be simplified by varying the query on this global and local scope for documents.

[Rz 4] The remaining part of the paper describes the proposed task and a theorized solution. First, we give summary of some related works. Subsequently, we describe our proposal, the methodology as well as evaluation.

¹ Eurovoc is available at <http://eurovoc.europa.eu>.

2. Background and Related Work

[Rz 5] Semantic Annotation formalizes and structures document with well-defined semantics specifically linked to a defined ontology [3]. Generally, annotation can aid structured organization of documents for optimized search. For instance, users may search information by well-defined general concepts that describe the domain of information need rather than use keywords.

[Rz 6] Ontology is a formal conceptualization of the world, capturing consensual knowledge [4]. It lists the concepts along with their properties and the relationships that exist between them. An example of a common knowledge resource used in NLP is Wordnet², a domain independent knowledge base of over 100,000 concepts in English in which a synset correspond to concept. Relationship between concepts are also well defined e.g. synonymy, antonymy, hyponymy etc. This study uses the Eurovoc descriptors as concept list. A concept can be annotated using lexicon based annotation or named entity identification. In the former, a dictionary of terms linked to each ontology class is maintained, reducing the annotation task to term matching. The latter uses recognition of named entities, with the entities mapped to the concept that gives their meaning.

[Rz 7] Semantic annotation can also be viewed as a classification task in which features are defined for each class and a Machine Learning (ML) classifier is trained to learn and group input document into its category [5][6][7]. Several SA systems have been implemented, for instance GoNTogle [8] uses weighted k Nearest Neighbor (kNN) classifier for document annotation and retrieval. A system widely used in semantic web domain is KIM [3] which assigns semantic descriptions to named entities (NE) in a text. The system is able to annotate and create hyperlinks to NEs inside a text and can then index and retrieve documents using these entities. Regular Expression (RE) has also been used to identify semantic elements in a text [9][10]. It works by mapping part of a text related to semantic context and matching the subsequent sequence of characters to create an instance of the concept. Application of these systems includes document retrieval especially in the semantic web domain [11][12]. Eneldo and Johannes [6] performed semantic annotation on legal documents for document categorization. Using Eurovoc concept descriptors on EurLex- a large database of legal documents, a ML classifier was trained for multi-label classification. While their work looks similar to our proposal, there are significant differences. First, their goal was a document categorization task considering and grouping a document as a whole while ours is two-fold, automatic segmentation of legal text into semantically coherent block and conceptual annotation of different segments with its rightful concept(s). Therefore, the proposed system better leans toward IR task than document categorization.

3. Research Motivation

[Rz 8] We hypothesize that conceptual zoning of document into semantically coherent blocks can enhance fine-grained IR and specifically enrich document retrieval systems as in eDiscovery procedures. Inspired by the recent successes in the area of Computational Linguistics (CL), we propose an automatic segmentation and semantic labeling of segmented portions of text. We introduce the idea of semantic scopes showing how important a concept is to a document, with the assumption that such semantic scopes can greatly enrich conceptual querying of big document corpus for IR,

² Wordnet is Available at <https://wordnet.princeton.edu/>.

once the documents in the corpus are well annotated. The idea can also help in scanning a voluminous legal text for specific part that is of utmost interest to the reader. For instance, in an IR task, the facts needed by a user in a document lie in a small portion of the whole document. Even if such facts are contained within a paragraph, a reader looking for specific information would have to read through the whole text in search of the fact, thus pilfering through «un-needed» information. With our idea, the portion of a document can be labeled according to its related concept; then the user (or further computational tools) have to only look for the text blocks that is tagged with a specific concept, this greatly speed up information gathering and simplify further, the task of information extraction from such processed text.

[Rz 9] To summarize the whole proposal without exploring the technical details, the system aims at achieving these stated goals. First, segmentation of text into semantically coherent blocks³ is done. We take idea from TextTiling [13][14] which divides text into contiguous, non-overlapping discourse units that correspond to the pattern of subtopics in a text; we advance this approach by incorporating some distributional analysis-informed heuristics (concept zoning⁴) to achieve the task.

[Rz 10] Secondly, we extract from a document the text portion(s) which best fit a specific information request based on concept filtering. To achieve this, we motivate the idea of conceptual scope of a specific document. For instance, a document can be described in terms of its *global* or *local* semantic scope; also, a concept could be local to a document or global to that document. We assume that a global scope concept has a higher argumentative weight in terms of its representation in the document while a local scope concept is lightly referred in terms of its representation in the document. The local or global scope representation further allow some degree of variableness in terms of how the system might be queried for information filtering and retrieval, providing a form of advanced search by enabling users to try different queries in respect of different context or condition and get responsive result sets. For instance, we may be able to vary the search query in terms of a global scope and one or more local scope(s) and get different result based on the local scope concept. As an example, a simple search query like «find all documents that talk about x in the context of y and not in the context of z » becomes possible, where x is a global scope category and y and z are local scope categories. As an example, let us assume that a document refer to three concepts «public-health», «Animal-feed» and «European-standard». We can also assume that public health has global scope while Animal feed and European standard both have local scope. Using this information, a user might attempt different queries as stated below and get different result:

- Retrieve the document(s) (and their specific parts) where the theme is generally about public health, while also talking about European standard
- Retrieve the documents and their parts where the general topic is about public health in relation to animal feed
- Retrieve the specific part(s) of a document that talks about European standards on animal feed in relation to public health.

³ Throughout the paper, we interchangeably use the words segment and block to mean the same thing. Also, concept and class are used interchangeably.

⁴ We take idea of zoning from the work of TEUFEL, S. (1999). Argumentative Zoning: Information Extraction from Scientific Text, which divides scientific papers into different sections called zones.

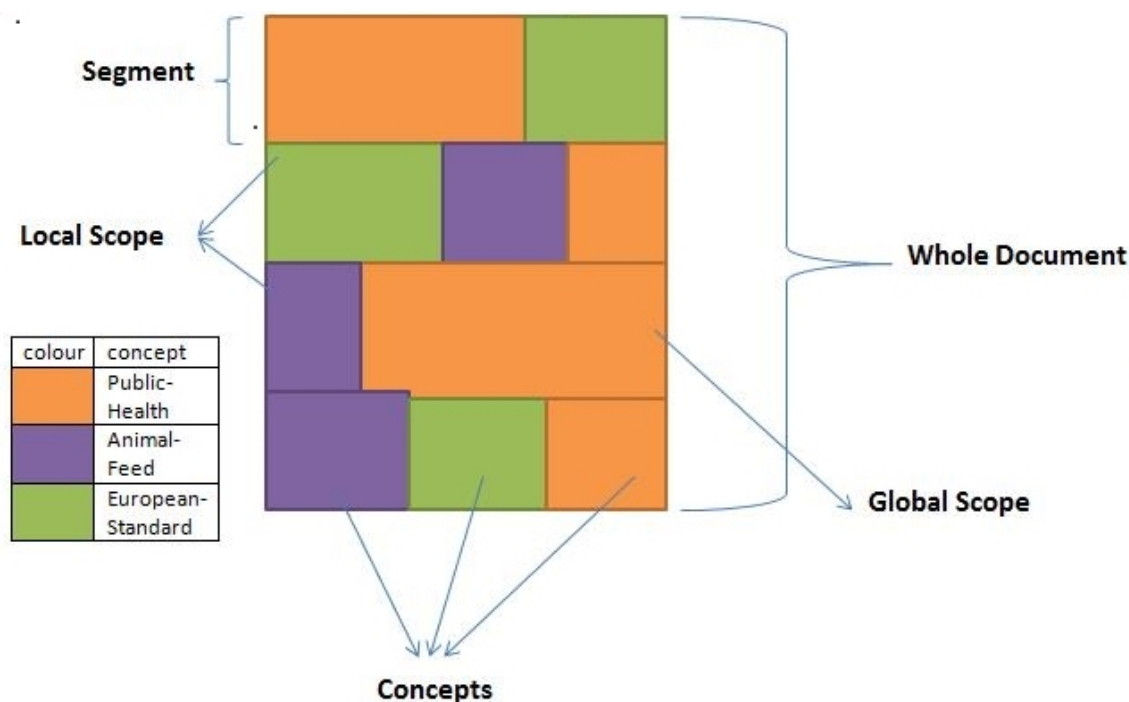


Fig. 1. Concept overlap and Semantic Scoping of Document

[Rz 11] The fig. 1 above gives a description of the semantic scoping of document into global and local scopes in terms of each concept’s representation of the document. Relying on a simple statistical procedure as visualized above, we can conclude that *public-health* is the most significant concept since it appears in every blocks. Also, its strength of representation in each of the blocks it appears is not negligible and thus, it can be labeled as the global concept. Whereas *Animal-feed* and *European standards* also appear in three of the four blocks, their representation in the document as a whole is incomparable to the former. We can also statistically determine the most significant concept that is local to each block as this can enable «*structured in-text scanning*». Thus, we may for example conclude that *Animal-feed* is global to the fourth block, even though it is itself local to the entire document as shown in the figure.

[Rz 12] This makes information filtering very advanced, for instance if a document talks about *Animal feed* but not in the context of *European standards* then it becomes irrelevant for the third query above and it is not retrieved. Also, if for example we have two documents A and B, with A having concepts *Public-Health* (*X*), *Animal-Feed* (*Y*) and *European-Standard* (*Z*) while B contains only concepts *Public-Health* (*X*) and *Animal-Feed* (*Y*) without *European-Standard*. Then as a proof of concept, queries⁵ of the form «SELECT FROM corpus WHERE concept = X, Y AND NOT concept = Z» can make a distinction between these concept overlaps and retrieve only document A while leaving out document B, even though it contains two of the mentioned concepts.

[Rz 13] We take a block as an information unit in a document, defined by different levels of granularity (i.e., sentence, paragraph or section holding many paragraphs). Our goal is to assign a label to each specific block with the assigned label corresponding to a specific concept in the ontology. The following processes are carried out:

⁵ The querying methodology is just a proof of concept and not implemented in this work.

1. Create concept profiles by looking at frequent and discriminant words calculated over texts-to-concepts occurrences (using TF-IDF).
2. Parse the text segment, sentence by sentence, calculating the similarity of each sentence with the concept profiles. Taking ideas from existing research on text segmentation such as Text-Tiling [13] which was improved in [14] and topic modeling [15] as well as TextRank [20], we identify the following flags in the text: (a) when the text in a block starts talking about a concept, (b) when it stops talking about the concept, and (c) when it starts talking about a new concept. Reference [16] contains a detail review of text segmentation approaches for the reader's interest.
3. Perform concept association, which maps a contiguous text block to a semantic concept that literarily gives a summary of its content.
4. Identify the global and local scopes of the concepts in a document by analyzing how the related concept block overlaps.

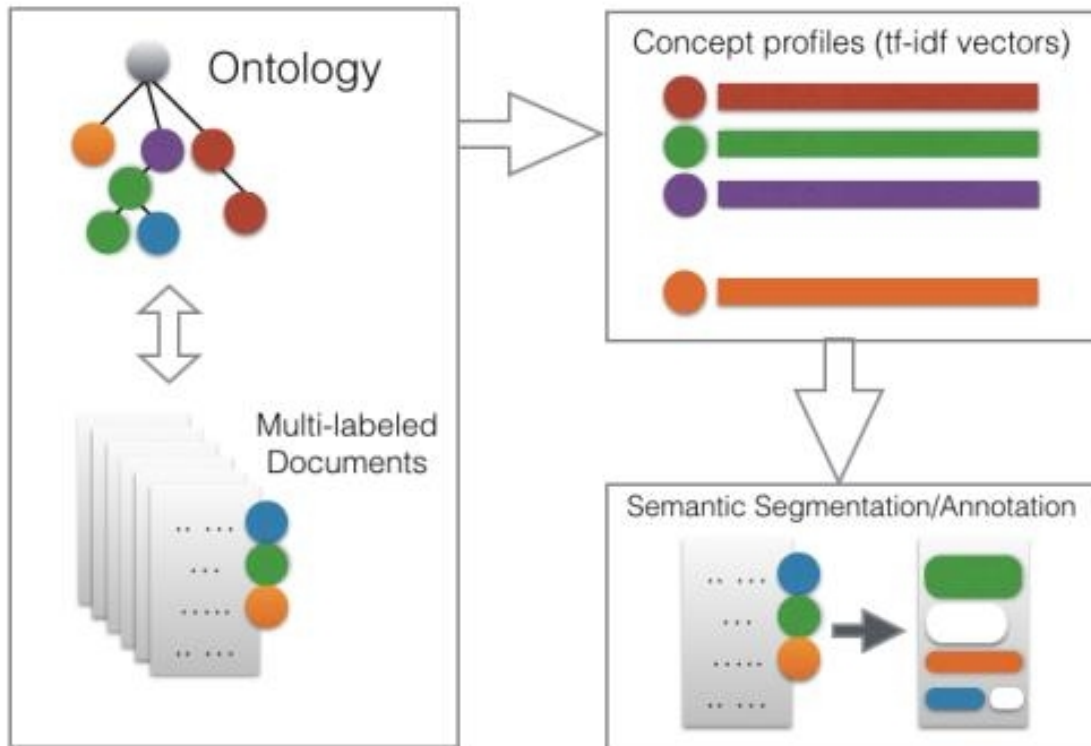


Fig. 2. Schematic representation of the proposed task

4. Methodology

[Rz 14] Let us consider a text t_i in a collection T of $n = |T|$ documents where $0 < i \leq n$. Each document is associated to a set of concepts taken from a thesaurus σ (i.e., a multi-labeled text collection).

[Rz 15] Thus we have some concepts $\{C_1, C_2, C_3, \dots, C_n\} \in \sigma$. We can formalize the task as that of approximating a target function

$$\Phi: B_{n-1, d} \times C\lambda \rightarrow \{1, 0\}$$

[Rz 16] that describes the conceptual mapping of a text segment. Where $B_{n-1, d}$ is a block or segment of text in a document d and n is a unique incremental number assigned as the identifier for each block. $C\lambda$ is a set of concepts according to a specified ontology. The goal is to find

$$\Phi (B_{n-1, d} \times C\lambda) = 1 \quad (1)$$

[Rz 17] Such that $C\lambda_i \equiv B_{i,d}$, that is, we want to have a mapping of a block to a concept successfully such that each block is associated to its concept⁶. The equivalence relation above ensures that a concept is attached to one or more blocks of a document d . It is possible for text blocks to share same concept as well as a single block being labeled by more than one concept.

[Rz 18] The first step is the creation of concept profiles, i.e., numeric vectors representing the contextual meaning of the concepts calculated through a TF-IDF weighting scheme over the concept-term matrix (which is built on the basis of the input multi-labeled document collection). This way, each concept C_p is associated to a vector v_p where $0 < p \leq |\sigma|$. Then, considering each text t_i as a sequence, $Seq_i = \langle s_1, s_2, \dots, s_k \rangle$ of sentences s , the idea is to label each sentence with zero or one element belonging to the set of concepts σ .

[Rz 19] This automatic segmentation/annotation of a text t_i is done as follows:

- Parse the sequence of sentences Seq_i , then

[Rz 20] For each S_j in Seq_i , calculate the cosine similarity between the frequency vector of S_j and each concept vector as below:

$$Sem_{dist}(C_p, S_j) = \frac{\vec{a}\vec{b}}{\|\vec{a}\|\|\vec{b}\|} \quad (2)$$

[Rz 21] Where

$$\vec{a}\vec{b} = \|\vec{a}\|\|\vec{b}\|\cos\theta \quad (3)$$

$$\cos\theta = \frac{\vec{a}\vec{b}}{\|\vec{a}\|\|\vec{b}\|} \quad (4)$$

[Rz 22] Such that \vec{a} represents the vector of concept C_p with \vec{b} representing the vectors of S_j in the text block.

[Rz 23] Then the similarity is obtained by the formula:

$$Sim(C_p, S_j) = \frac{1}{Sem_{dist}(C_p, S_j)} \quad (5)$$

[Rz 24]

- If the similarity between a concept c_p and a sentence s_j is z times higher⁷ than the rest of concept-to-sentence similarities, then s_j is associated to c_p , otherwise the sentence is not associated to any concept. The parameter z is a predetermined threshold value strictly for

⁶ i signifies an iterative number, incrementing over the sets of concepts and blocks in a document d .

⁷ Unique parameter of the method

decision making. While this value can be varied, it is set to 2 by default, making it possible for a sentence to be mapped to its most similar concept.

- Finally, semantically-contiguous sentences (i.e., sentences which are contiguous and associated to a single concept) will represent semantically-coherent segments, which are the final result of the in-text concept annotation task.
- In case the method returns an empty set of segments or an incomplete coverage of the concepts associated to the text t_i , it restarts by step 1 with $z = z/2$.

[Rz 25] A schematic view of the task is shown in fig. 2 above.

4.1. Concept Analysis

[Rz 26] Concept descriptors could range from unigram, bigrams to n-grams. Similar to query expansion, if a concept descriptor⁸ has more than one word, we break the n-gram terms into the constituent words in a process called lexical expansion. The goal of lexical expansion is to retrieve semantically similar words such as synonyms, relying on a knowledge base such as Wordnet. This implies that the document need not contain in explicit terms the constituent words that make up the concept as part of the document's key terms. We used a concept profile, which contains keywords from the deflated n-grams, as well as the synonyms of each of the words for each concept. Weights are assigned to each synonym based on their path distance to each of the lexically expanded terms of concept, leading to the selection of the best ranked synsets. We then perform a form of synonym merging on the ranked terms for each concept, which combines these terms by collapsing them into a single block of information unit. With the lexical expansion, vectorization [17] [18] is done to build a vector of terms that describe the information content of each concept using.

[Rz 27] For simplicity, if σ is a set such that $\sigma = \{c_1, c_2, \dots, c_k\}$ which is a list of all the concepts for that document. A concept c_i may also have multiple terms e.g. Public-health, each of this, along with its synonyms is a term in the vector space whose frequency of occurrence in each text block is quantified. TF-IDF is used to create vector representations of each concept as well as each text block. Each component of a vector corresponds to the TF-IDF value of a particular term in the text block dictionary. Dictionary terms that do not occur in a block are weighted zero, taking the representation as query vectors that can be compared to vectors of documents (here each text block). The semantic distance between a concept and the text block is calculated using equation 2 and the similarity between a concept and the text bloc is calculated according to equation 5. Iteration is made over each concept in σ , calculating how similar it is to the text block and if similar based on a fixed parameter z , such concept is tagged with that block. The process is repeated for all the blocks in the document which leads to the idea of concept(s)-block tagging. Fig. 3 below shows the general system architecture.

⁸ Take for instance public-health which can be dissolved into public and health.

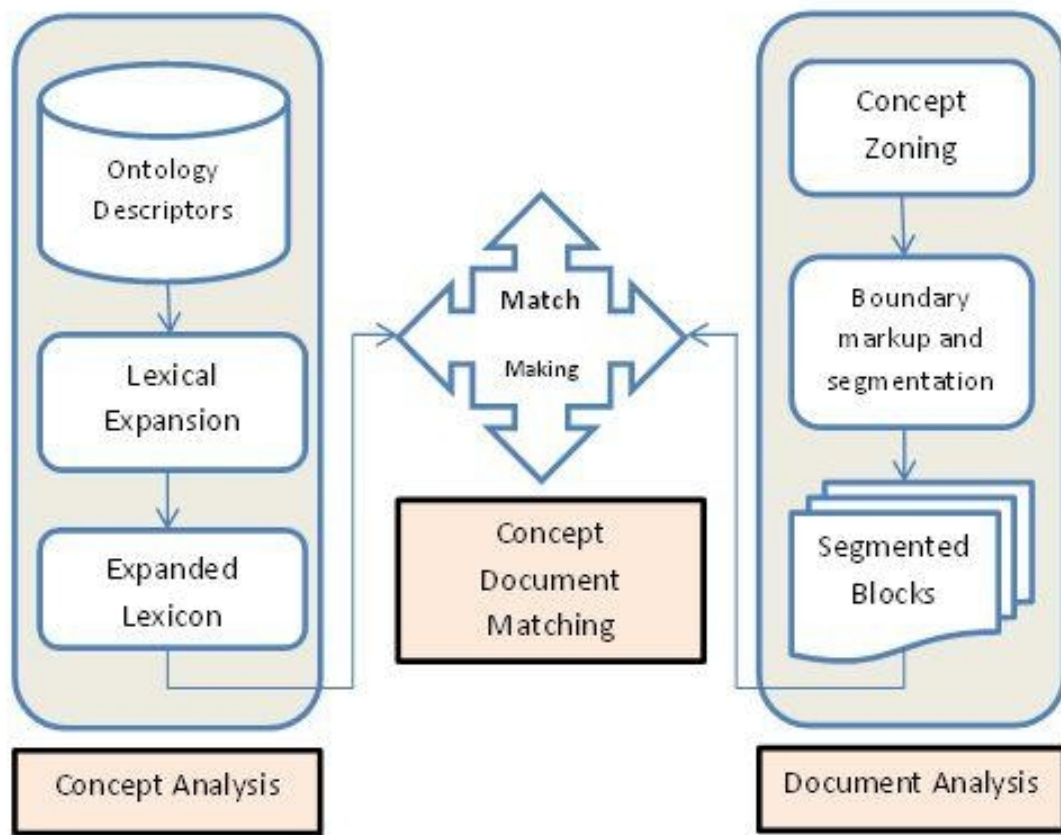


Fig. 3. Proposed System Architecture

4.2. Concept Zoning for Text Segmentation

[Rz 28] Given a text t_i and its sequence of sentences $Seq_i = \langle s_1, s_2, \dots, s_k \rangle$. We perform concept zoning, which aggregates all the sentences or paragraphs of the text that semantically align into a group. This semantically aligned group is taken as a block, each of which can be directly mapped to the concept descriptor from an ontology. To segment text into coherent semantic units, we employed word clustering [19]. We used the popular K-means clustering which is based on Lloyd's algorithm [21]. Here, clusters are formed of document parts that appear to be highly similar. First, text is converted into a bag-of-words with each sentence extracted as a micro-document. K-means algorithm is then applied to cluster related words together. Since the input is sentences, K-Means clusters similar sentences. A requirement of K-Means algorithm is the number of clusters to be specified in advance [22]. We here assumed a fixed cluster size of 3. In theory, the single cluster containing the entire collection is represented by the document itself and serves as the root of the tree while the coalesced segments are leaves in the bigger tree. Each cluster is a text segment representing a semantically coherent unit of the document. The k value can be varied in order to increase the number of clusters and of course, document segments.

4.3. Concept-Document Mapping

[Rz 29] Matching a given concept to a block in a text is reduced to a simple semantic relatedness task between the term-document vector of the expanded lexicon of each concept and the bag-of-word vector of each segment. TF-IDF weighing scheme was used while cosine similarity was used to measure the semantic distance between the vectors as explained earlier using the formula

$$Sim(C_p, S_j) = \frac{1}{Sem_{dist}(C_p, S_j)} \quad (5)$$

[Rz 30] For each of the segment, the system iterates over all the term-document vectors of the expanded lexicon of each concept, measuring the distance. Similar vectors imply some level of relatedness between the concept and the segment and such segment is tagged with the concept.

5. Evaluation

[Rz 31] We randomly sampled 5 documents from Eur-Lex database of legal documents. Eur-Lex⁹ is an open and regularly updated database of over 3 million EU legal documents, covering EU treaties, regulations, legislative proposals, case-law, international agreements, EFTA documents and some other public documents of interest to EU operations. For each of the documents in the EurLex database, a list of concept(s) describing the document is already listed. We used Eurovoc concept descriptors as ontology. Eurovoc is a multilingual and multidisciplinary thesaurus. Most language versions contain over 6883 preferred concept descriptors and up to 10,592 non-preferred concept terms, organized hierarchically into 21 domains that is of interest to EU's parliaments. Currently, it is available in 26 European languages. We evaluated the system on the task of conceptual tagging.

- Conceptual Tagging: This task measures the performance of the system in correctly labeling a text segment with a concept. EurLex documents are pre-classified with some concepts. We required a volunteer to identify and manually annotate portions of the text that talks about each concept classified for each document. We measured the performance of the system against annotations from human judgment and got an accuracy of 62%.

6. Conclusion

[Rz 32] We have described in this paper, a work-in-progress on a task that involves semantic in-text labeling of text blocks annotated with ontology concepts.

[Rz 33] The task employs the use of Eurovoc ontology, a multilingual thesaurus continuously updated by the EU publication's office. The task considers a new approach aimed at enhancing information filtering within text by advancing the classification tasks with semantic annotation. We proposed a formalization of the task and a baseline approach to solve it.

[Rz 34] The task has a lot of potential applications in IR, text segmentation, topic modeling and text summarization as well as argumentation mining. For instance, a user who is more interested

⁹ Available at <http://eur-lex.europa.eu/content/welcome/about.html>.

in a particular information in a document can easily specify such information need through an concept that describes such need and the system is able to extract the specific portion containing the requested information need. This is made possible with semantic tag associated with text portion(s), showing their semantic alignment to each ontology terms. Thus, users need not pilfer through unwanted information. We have proposed an evaluation of the system on conceptual tagging of text, benchmarking the system with manual annotations from human and using such manual annotations as Gold standard. The result obtained shows that the approach is promising but requires a bigger and thorough evaluations to be able to compare with results from existing systems. Subsequent works will provide a detailed experimental analysis and evaluation results of our approach in different context and domain, for instance, we will explore in deep, the information retrieval task aspect of our work as well text segmentation.

7. References

- [1] EDRM, Electronic Discovery Reference model, <http://www.edrm.net>.
- [2] The 2008 Socha-Gelbmann Electronic Survey Report (2008). <http://www.sochaconsulting.com/2008survey.php>.
- [3] POPOV, B., KIRYAKOV, A., KIRILOV, A., MANOV, D., OGNANOFF, D., & GORANOV, M. (2003). KIM—semantic annotation platform. In *The Semantic Web-ISWC 2003* (pp. 834–849). Springer Berlin Heidelberg.
- [4] KIYAVITSKAYA, N., ZENI, N., MICH, L., CORDY, J. R., & MYLOPOULOS, J. (2006). Text mining through semi automatic semantic annotation. In *Practical Aspects of Knowledge Management* (pp. 143–154). Springer Berlin Heidelberg.
- [5] ASOOJA, K., BORDEA, G., VULCU, G., O'BRIEN, L., ESPINOZA, A., ABI-LAHOUD, E., & BUTLER, T. (2014). Semantic Annotation of Finance Regulatory Text using Multilabel Classification.
- [6] DAELEMANS, W., & MORIK, K. (eds.). (2008). *Machine Learning and Knowledge Discovery in Databases: European Conference, Antwerp, Belgium, September 15–19, 2008, Proceedings* (Vol. 5212). Springer.
- [7] BUABUCHART, A., METCALF, K., CHARNESS, N., & MORGENSTERN, L. (2013). Classification of Regulatory Paragraphs by Discourse Structure, Reference Structure, and Regulation Type. In *JURIX* (pp. 59–62).
- [8] BIKAKIS, N., GIANNOPOULOS, G., DALAMAGAS, T., & SELLIS, T. (2010). Integrating keywords and semantics on document annotation and search. In *On the Move to Meaningful Internet Systems, OTM 2010* (pp. 921–938). Springer Berlin Heidelberg.
- [9] LACLAVÍK, M., CIGLAN, M., SELENG, M., & KRAJEI, S. (2007). Ontea: Semi-automatic pattern based text annotation empowered with information retrieval methods. *Tools for acquisition, organisation and presenting of information and knowledge: Proceedings in Informatics and Information Technologies, Kosice, Vydavateľstvo STU, Bratislava, part, 2* (pp. 119–129).
- [10] LACLAVIK, M., SELENG, M., GATIAL, E., BALOGH, Z., & HLUCHY, L. (2007). Ontology based text annotation-OnTeA. *Frontiers in Artificial Intelligence and Applications* (pp. 154, 311).
- [11] HANDSCHUH, S., & STAAB, S. (2002). Authoring and annotation of web pages in CREAM. In *Proceedings of the 11th international conference on World Wide Web* (pp. 462–473). ACM.

- [12] DILL, S., EIRON, N., GIBSON, D., GRUHL, D., GUHA, R., JHINGRAN, A., & ZIEN, J. Y. (2003). A case for automated large-scale semantic annotation. *Web Semantics: Science, Services and Agents on the World Wide Web* (pp. 115–132).
- [13] HEARST, M. A. (1993). *TextTiling: A quantitative approach to discourse segmentation*. Technical report, University of California, Berkeley, Sequoia.
- [14] HEARST, M. A. (1997). TextTiling: Segmenting text into multi-paragraph subtopic passages. In *Computational linguistics* 23, no. 1 (1997): 33–64.
- [15] RIEDL, M. & BIEMANN C. (2012). Text segmentation with topic models. In *Journal for Language Technology and Computational Linguistics* 27, no. 1 (2012): 47–69.
- [16] LLORET, E. (2009). Topic Detection and Segmentation in Automatic Text Summarization. In *Focus Journal*.
- [17] CLARK, S. (2014). Vector space models of lexical meaning (to appear). In *Handbook of Contemporary Semantics*. Wiley-Blackwell, Oxford.
- [18] TURNEY, P. D. & PANTEL, P. (2010). From Frequency to Meaning: Vector Space Models of Semantics. In *Journal of Artificial Intelligence Research* 37 (pp. 141–188).
- [19] PANTEL P., LIN D. (2002). Discovering word senses from text. In *Proc 8th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 613–619). Edmoton, Canada.
- [20] MIHALCEA, R., & TARAU, P. (2004). TextRank: Bringing order into texts. Association for Computational Linguistics.
- [21] HARTIGAN, J. A., & WONG, M. A. (1979). Algorithm AS 136: A k-means clustering algorithm. In *Applied statistics* (1979): 100–108.
- [22] KANUNGO, T., MOUNT, D. M., NETANYAHU, N. S., PIATKO, C. D., SILVERMAN, R. & WU, A. Y. (2002). An efficient k-means clustering algorithm: Analysis and implementation. In *Pattern Analysis and Machine Intelligence, IEEE Transactions on* 24, no. 7 (2002): 881–892.