

Crime Knowledge Extraction: an Ontology-driven Approach for Detecting Abstract Terms in Case Law Decisions

Silvana Castano
silvana.castano@unimi.it
Università degli Studi di Milano
Computer Science Department
Milan, Italy

Alfio Ferrara
alfio.ferrara@unimi.it
Università degli Studi di Milano
Computer Science Department
Milan, Italy

Mattia Falduti
mattia.falduti@unimi.it
Università degli Studi di Milano
Computer Science Department
Milan, Italy

Stefano Montanelli
stefano.montanelli@unimi.it
Università degli Studi di Milano
Computer Science Department
Milan, Italy

ABSTRACT

In this paper, we present *CRIKE*, a data-science approach to automatically detect concrete applications of legal abstract terms in case-law decisions. To this purpose, *CRIKE* relies on the use of the *LATO ontology* where legal abstract terms are properly formalized as concepts and relations among concepts. Using *LATO*, *CRIKE* aims at discovering how and where legal abstract terms are applied by judges in their legal argumentation. Moreover, we detect the terminology used in the text of case-law decisions to characterize concrete abstract-term instances. A case-study on a case-law decisions dataset provided by the Court of Milan, Italy, is also discussed.

CCS CONCEPTS

• Information systems → Ontologies.

KEYWORDS

legal ontology, legal-term extraction, case-law analysis

ACM Reference Format:

Silvana Castano, Mattia Falduti, Alfio Ferrara, and Stefano Montanelli. 2019. Crime Knowledge Extraction: an Ontology-driven Approach for Detecting Abstract Terms in Case Law Decisions. In *Seventeenth International Conference on Artificial Intelligence and Law (ICAIL '19)*, June 17–21, 2019, Montreal, QC, Canada. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/3322640.3326730>

1 INTRODUCTION

Law is general and abstract by definition. On the opposite, court case law decisions are specific and concrete, in that they provide a peculiar interpretation of law applied to the considered single

cases. Legal interpreters, such as for example judges and lawyers, are daily involved in analysis and evaluation of court case law with the aim to extract/derive possible suggestions for incoming case applications by relying on the experience of past applications that can be considered as a sort of consolidated *legal knowledge*.

According to the Italian law, the legal terminology can be distinguished into three main categories, that are i) *statutory terms*, i.e., terms directly or indirectly defined by law; examples of statutory terms are public officer, illicit drug, and consumer; ii) *descriptive terms*, i.e., terms featuring actions, human activities, and any real-life object; examples of descriptive terms are escape, car, and year; iii) *abstract terms*, i.e., terms featuring something indeterminate that requires a concrete application for being really defined; examples of abstract terms are good faith, long-term cohabitation, and dangerous driving. Consider the abstract schema of a legal action provided in Figure 1. When a new case law is received for judgement, the

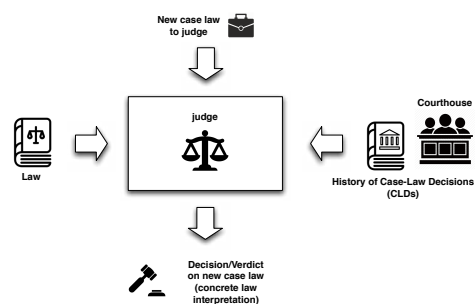


Figure 1: The abstract schema of a legal action

expected evaluation process has to take into account i) the law, for understanding the terms, either statutory, descriptive, or abstract, that can be relevant for the current case, and ii) the history of case-law decisions, for detecting possible relevant interpretations and concrete applications of abstract terms that can be useful to support the decision/verdict to eventually deliver.

In this paper, we present *CRIKE* (CRIME Knowledge Extraction), a data-science approach to detect concrete applications of legal

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ICAIL '19, June 17–21, 2019, Montreal, QC, Canada

© 2019 Association for Computing Machinery.

ACM ISBN 978-1-4503-6754-7/19/06...\$15.00

<https://doi.org/10.1145/3322640.3326730>

abstract terms in large case-law decisions. To this purpose, CRIKE relies on the use of LATO (Legal Abstract Term Ontology) where legal terms are properly formalized as concepts and relations among concepts. Using LATO, CRIKE aims at discovering how and where legal abstract terms are applied by judges in their legal argumentation. A case-study provided by the Court of Milan, Italy, is finally discussed to show the results of the CRIKE application on a real dataset of 207 case-law decisions, where documents are different in terms of redaction method and legal content and a manual annotation step with the support of a domain expert has been applied.

The paper is organized as follows. In Section 2, the CRIKE approach is introduced. The LATO ontology and the CRIKE techniques for legal knowledge extraction are discussed in Section 3 and 4, respectively. In Section 5, CRIKE support to practices and preliminary experimental results are presented. Related work are discussed in Section 6. Concluding remarks are provided in Section 7.

2 THE CRIKE APPROACH

The CRIKE approach (see Figure 2) is conceived to support extraction of legal knowledge from a (possibly large) dataset of Case-Law Decisions (CLDs) coming from different, official sources, such as for example First Grade and Court of Appeal judgements. CRIKE embeds the LATO ontology where relevant law concepts of a given domain of interest are properly formalized. To enforce knowledge extraction, CRIKE exploits a given dataset of CLDs in input by adopting a conventional data-science process where each CLD is indexed and stored in a digital format. In particular, the CLDs of our dataset are acquired from the Court and the Court of Appeal of Milan and they are usually provided in image format with highly heterogeneous quality. The indexing and storage activity exploits data cleaning and tokenization techniques to obtain a pure textual version of each CLD as well as a focused set of metadata. By exploit-

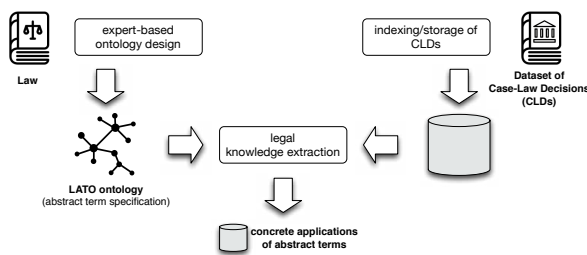


Figure 2: The CRIKE approach

ing the indexed CLDs metadata, knowledge extraction is enforced with the aim at classifying a CLD with respect to the LATO ontology knowledge. In particular, extraction is focused on detecting the concrete applications of legal abstract terms in the text of the considered CLDs. The crucial idea of CRIKE is that the detection of a given abstract term *AT* is not only concerned with the recognition of single terms featuring *AT*, but also with the recognition of terms associated with the ancillary concepts related to *AT*, that we call *abstract-term context*.

Motivating example. Consider the Italian law about drugs and related drug offenses, as reported in [10]. According to the Italian criminal order, *the Consolidated Law, adopted by Presidential Decree No 309 on 9 October 1990 and subsequently amended, provides the legal framework for trade, treatment and prevention, and prohibition and punishment of illegal activities in the field of drugs and psychoactive substances. Drug use in itself is not mentioned as an offense. [...] The threshold between personal possession and trafficking is determined by the circumstances of the specific case (e.g., the act, possession of tools for packaging, different types of drug possessed, number of doses in excess of average daily use, means of organization). The penalty for supply-related offenses, such as production, sale, transport, distribution or acquisition, depends on the type of drug. However, when the offenses are considered minor because of the means, modalities or circumstances, the terms of imprisonment are lower. Evaluating whether or not the offense is minor should take into account a set of “ancillary” elements such as the mode of action, possible criminal motives, quality and quantity of drug possessed, the character of the offender, conduct during or subsequent to the offense, and the family and social conditions of the offender.* The notion of **minor offense** is an example of abstract term in the above law quotation. A precise definition of circumstances and related threshold quantities to associate with the notion of minor offense is not available/possible in the (abstract) law. Given a specific criminal charge of drug possession, the final decision/verdict is based on the specific interpretation of the abstract term “minor offense” where the specific circumstances and quantities of the considered case represent a concrete application of the corresponding abstract term.

3 LEGAL KNOWLEDGE REPRESENTATION

In order to formalize the knowledge related to abstract terms and their interpretation, we introduce LATO in CRIKE. LATO is a legal ontology where relevant law terms to exploit for CLDs analysis and knowledge extraction are defined; it contains concepts to represent general law terms, either abstract, statutory and descriptive terms¹.

LATO is manually defined by domain experts and implemented according to the SKOS formalism. In particular, the concept hierarchy is based on a root concept Term with three main subconcepts, namely AbstractTerm, DescriptiveTerm, and StatutoryTerm (see Figure 3(a)). In addition to general law terms, the LATO ontology contains concepts that represent the Italian legislative structure, such as for example the concepts Law, LawArticle, and LawParagraph. Furthermore, the concepts Conviction and Discharge are also specified in LATO to represent the possible Court decisions (i.e., the verdict) of a given case law. In particular, the concept Conviction denotes a verdict in which the Court judges the defendant guilty, while the concept Discharge denotes a verdict in which the facts have a penalty relevance, but no punishment is finally delivered. Finally, the concepts Quantity and UnitOfMeasure are defined in LATO for allowing to represent the quantitative estimation of substances that can appear in legal documents.

AbstractTerm is the core concept of the LATO ontology since it represents the target of the knowledge extraction functionalities of CRIKE. The related construct of SKOS is exploited to enrich the

¹In the current version, the LATO ontology covers the main terms of the Italian criminal code.

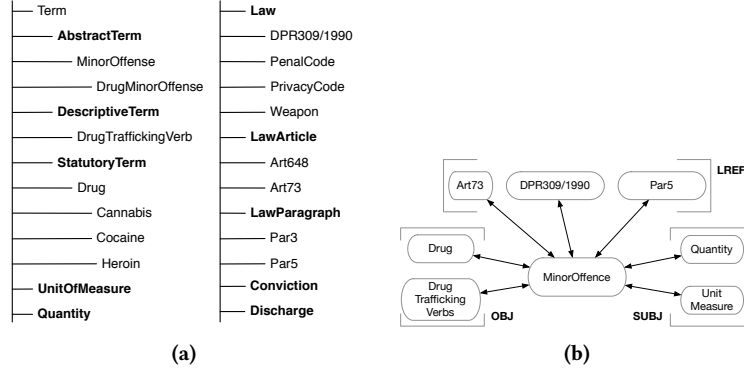


Figure 3: (a) Excerpt of the LATO concept hierarchy; (b) Example of concept definition for the abstract term MinorOffense

specification of an abstract term AT by formalizing the ontology relationships between AT and the other concepts of the LATO ontology composing its context. In particular, given a considered abstract term AT , related is used to connect AT to ancillary concepts of LATO representing i) an *objective judgment element* OBJ usually expressed through the connection of AT with a descriptive/statutory concept; ii) a *subjective quantitative evaluation* $SUBJ$ usually expressed through a relationship between AT and Quantity/UnitOfMeasure concepts; and iii) a *legislative reference* $LREF$ usually denoted with a connection of AT with a specific law or regulation (i.e., Law, LawArticle, and LawParagraph concepts).

According to SKOS, each LATO concept is associated with a *preferred label* ($prefLabel$) as well as with one or more *alternative labels* ($altLabel$) and *hidden labels* ($hiddenLabel$) to enrich the concept definition with a label-set of literal descriptions that is very useful for subsequent knowledge extraction, to capture possible synonyms, acronyms, and abbreviations in the text of CLDs.

Example. An example of SKOS definition for the abstract term $AT = MinorOffense$ is shown in Figure 3(b) according to the Italian drug-trafficking law. $MinorOffense$ is related to the concepts Drug and DrugTraffickingVerb that represent the OBJ relationships since they are subconcepts of StatutoryTerm and DescriptiveTerm, respectively. The relationships with the concepts Quantity and UnitOfMeasure represent the subjective judge evaluations $SUBJ$. The concepts Par5, Art73, and DPR309/1990 are subconcepts of the LawParagraph, LawArticle, and Law, respectively, and they express the legal references $LREF$ of $MinorOffense$ in the Italian criminal code where the drug trafficking crime is defined.

4 KNOWLEDGE EXTRACTION IN CRIKE

Knowledge extraction in CRIKE is based on the idea to exploit text analysis techniques for detecting the concrete applications of legal abstract terms belonging to LATO throughout the stored/indexed case-law decisions CLDs. To this end, for a given abstract term AT , we introduce the notion of *abstract-term context* Ctx_{AT} containing, besides the AT term, all the concepts of LATO that are ancillary to

AT , namely OBJ , $SUBJ$, or $LREF$ concepts:

$$Ctx_{AT} = \{C_i \mid r(AT, C_i)\}$$

where $r(AT, C_i)$ denotes a SKOS related relationship between the abstract term AT and the concept C_i .

For each concept $C \in Ctx_{AT}$, we define the *concept label set* L_C that contains the whole set of labels, either preferred, alternative, or hidden, associated with C . Furthermore, based on the notion of L_C , we define the *extended label set* \mathcal{L}_C where the concept label set of C is enriched by including the concept label set of the concepts C_j subsumed by C :

$$\mathcal{L}_C = L_C \cup \{L_{C_j} \mid C_j \subseteq C\}$$

Consider the goal to detect the concrete applications of a certain abstract term AT in a dataset of case-law decisions $CLDs$. CRIKE knowledge extraction is enforced by exploiting the extended label sets \mathcal{L}_C of the concepts in the context Ctx_{AT} . For each document $d \in CLDs$, we define a vector representation \vec{d} where each element corresponds to a concept in the context Ctx_{AT} . The value $d[i] \in \vec{d}$ is set to 1 when a *label hit* is detected, meaning that at least one occurrence of a label in \mathcal{L}_{C_i} is found in d for the concept $C_i \in Ctx_{AT}$, and 0 otherwise (i.e., *label miss*). A threshold based mechanism is defined to specify the minimum number of label hits required to consider that a concrete application of the abstract term AT is detected in the document d .

Example. Consider the abstract term $AT = MinorOffense$ and the corresponding context $Ctx_{MinorOffense} = \{Drug, DrugTraffickingVerb, DPR309/1990, Art73, Par5, Quantity, UnitOfMeasure\}$. Moreover, consider the extended label set $\mathcal{L}_{Drug} = L_{Drug} \cup \{L_{Cocaine}, L_{Heroin}, L_{Cannabis}\}$. In Figure 4, we show an example of knowledge extraction based on the concepts and corresponding extended label sets in the context $Ctx_{MinorOffense}$. An example of vector-based document representation for the abstract term $AT = MinorOffense$ is shown in Figure 5. If we consider a threshold of 80% of label hits, we have that a concrete application of $MinorOffense$ is detected in document d_1 since 6 hits are found over the available 7 concepts in $Ctx_{MinorOffense}$.

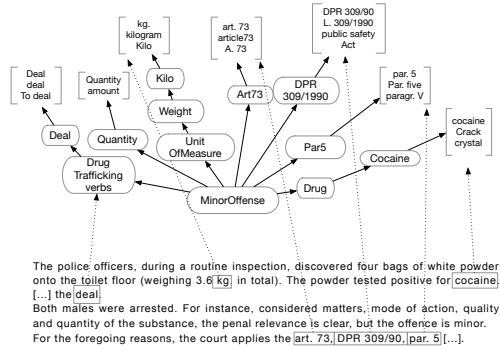


Figure 4: Example of knowledge extraction for the abstract term MinorOffense

	Drug	DrugTraffickingVerb	DPR309/90	Art73	Par5	Quantity	UnitOfMeasure
d_1	1	0	1	1	1	1	1
d_2	1	1	1	0	0	1	0

Figure 5: Example of label hits for the abstract term MinorOffense

5 CRIKE SUPPORT TO PRACTICES AND PRELIMINARY EXPERIMENTAL RESULTS

In the following, we introduce possible applicative issues for the proposed CRIKE approach and we discuss some preliminary results obtained on a dataset of Italian case-law decisions.

5.1 Applicative issues of CRIKE

The CRIKE approach can be exploited to support the judge activities by providing knowledge automatically extracted from consolidated case-law decisions. In particular, we envisage the following applicative scenarios.

Knowledge-assisted verdict writing. CRIKE can be exploited to support the judge in the preparation of new case-law decisions. The idea is that the LATO ontology allows to classify the available case-law decisions, so that a similarity-based retrieval process can be enforced by considering both the contents of the new incoming case and the ontology knowledge extracted from previous decisions. This way, a judge is supported in writing by receiving hints and suggestions from law interpretations and concrete law applications on similar cases.

History-based verdict prediction. The knowledge extracted by CRIKE can be exploited to train a machine learning approach, with the aim to predict the possible decision on a new incoming case-law to judge according to the case context (e.g., circumstances, quantities) in relation with features and decisions of previous case-law decisions. Such a prediction can be considered as another form of judge support in verdict preparation based on statistical evidence instead of similarity.

Legal analytics. CRIKE allows to enforce analytics over the knowledge extracted from case-law decisions with the aim to detect possible trends and common law interpretations in presence of

	Minor offence		Conviction			Discharge			
True	0	103	2	0	32	15	0	187	0
	1	12	90	1	3	157	1	2	18
		0	1		0	1		0	1
		Predicted			Predicted			Predicted	

Figure 6: Results of CRIKE knowledge extraction for the concepts of minor offence, conviction, and discharge

certain context features. Due to the nature of considered data, it is important to note that ethical issues are particularly important when developing techniques for analytics over legal data, so that fairness is preserved and bias/discriminations are properly avoided.

5.2 Experimental results

For experiments on the CRIKE approach, we collected a dataset of 226.413 decision documents provided by the Courthouse of Milan, Italy. The dataset contains either decision documents about the First Grade of decision (i.e., 123.186 documents) and Court of Appeal decisions (i.e., 103.227 documents). Each decision document is associated with a unique identification number and it is provided as a PDF image with possible handwritten annotations, stamps, and symbols/abbreviations. Documents are heterogeneous in terms of both formatting aspects (e.g., fonts and page layout) and editorial content (e.g., document organization and outline). Due to privacy regulations and public authority orders, the dataset is not public and it cannot be shared.

For a preliminary evaluation of the knowledge extraction functionalities of CRIKE, we consider a subset of 207 decisions about drug trafficking and minor offense with related verdicts. The 207 documents have been manually annotated by a domain expert with the aim to distinguish where i) a minor offense has been recognized by the judge, and ii) a certain verdict has been delivered, either conviction (*Conv*) or discharge (*Disc*). In particular, given 207 decisions, the expert annotated 102 documents where minor offenses were recognized, 160 documents had a conviction verdict, and 20 has a discharge verdict (the remaining 27 decisions has been annotated with an acquittal verdict). The expert annotations has been exploited as ground truth for evaluation against the CRIKE results obtained by relying on the specifications of the abstract concepts MinorOffense, Conviction, and Discharge contained in the LATO ontology.

The results of knowledge extraction performed by CRIKE over the considered dataset 207 decision documents is shown in Figure 6.

In the results, we compare the document classification based on the ground truth (True side of the confusion matrix) and the CRIKE classification (Predicted side). In the experiment, the CRIKE results are promising. Indeed, we note that only 14 documents are mis-recognized by CRIKE for the concept minor offense (6.8% missed). Similar performance are obtained for the concept conviction (18 missed, 8.7%). Very promising results are obtained for the concept discharge where only 2 documents are mis-recognized, probably due to the very specific and recurring terminology associated with the discharge concept in the documents.

As a further result, in Figure 7, we show the distribution of concept labels in the context $Ctx_{MinorOffence}$ that has been recognized in the documents classified as minor offense. This way, it

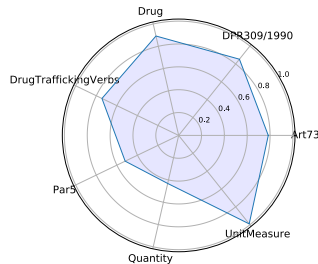


Figure 7: The concept labels recognized in documents about minor offense

is possible to have an analytic insight about which labels of the abstract concept minor offense are actually recognized in the considered documents. We note that CRIKE correctly recognizes the law references of the concept minor offense (i.e., Art73 and DPR309/1990) in more than the 80% of considered documents. Moreover, the concept labels Drug (including DrugTraffickingVerbs) and UnitMeasure are recognized in almost all the documents. On the opposite, Par5 and Quantity are correctly recognized in less than the 50% of documents. This result suggests that the context specifications of these concepts need to be better enriched in the LATO ontology.

6 RELATED WORK

Work related to the issues addressed in CRIKE regards legal argumentation mining and legal ontology design. Legal argumentation mining refers to the capability to automatically detect and classify the role of possible argumentative units within a considered legal case text [1]. In [9], authors propose to mine statutory texts by using natural language processing and supervised machine learning techniques. More recently, the LUIIMA approach has been proposed to focus on extraction of evidential reasoning from a court decision dataset [4]. Moreover, a particularly relevant contribution is provided in [8] about extraction of case law sentences for argumentation of statutory terms.

A survey on legal ontology design is presented in [1], where a special focus is given to representation of legal concepts in type systems. In [3], the notion of mutual consensus is introduced to support the specification of concepts and relations about contract formation. An application example based on a corpus of Italian legal texts is presented in [6], where the results of exploiting a learning system are provided. A further specification of a legal ontology using ONTOLINGUA is presented in [12]. Furthermore, in [11], authors present the LOIS project (Lexical Ontologies for Legal Information Sharing), and discuss a methodology for building a multilingual semantic lexicon for law able to be used both as a source of semantic metadata and as an external tool for cross lingual retrieval. On that topic, in [7], a methodology to automatically create an

OWL ontology from a set of legal documents is presented. In [2], an automated approach based on statistical analysis is described, for identification of core concepts and relations in a corpus of legal texts. Natural Language Processing (NLP) techniques are proposed in [5], to extract concepts and relations among legal concepts, with the aim to build an ontology for legal information retrieval.

Original contribution of the proposed CRIKE approach is related to the enforcement of a data-science process with the support of an expert-based law ontology to extract knowledge from CLDs. A further peculiar feature of CRIKE is related to the formalization of an abstract term as a legal ontology concept with a corresponding context of related concepts. Ontology concepts with associated contexts are used to drive the identification of concrete applications of corresponding abstract terms in the text of CLDs.

7 FUTURE WORK

In this paper, we presented the CRIKE approach for legal knowledge extraction.

Different research directions are currently being investigated. On the one side, we are working on a bootstrapping approach to enforce enrichment of the LATO ontology, so that the context of abstract terms can be progressively augmented with new relevant terms and literals as long as they are detected in CLDs during extraction. On the other side, machine learning techniques are being developed to enforce a supervised classification of CLDs based on abstract terms, by exploiting a training set of CLDs manually annotated by domain experts. A further research topic is related to the identification of ethical issues involved in the CRIKE data-science process for CLD analysis and classification.

REFERENCES

- [1] Kevin D Ashley. 2017. *Artificial Intelligence and Legal Analytics: New Tools for Law Practice in the Digital Age*. Cambridge University Press.
- [2] Enrico Francesconi, Simonetta Montemagni, Wim Peters, and Daniela Tiscornia. 2010. *Integrating a Bottom-Up and Top-Down Methodology for Building Semantic Resources for the Multilingual Legal Domain*. Vol. 6036. Springer, 95–121.
- [3] Anne Gardner. 1987. *An Artificial Intelligence Approach to Legal Reasoning*. MIT Press, Cambridge, MA, USA.
- [4] Matthias Grabmair, Kevin D Ashley, Ran Chen, Preethi Sureshkumar, Chen Wang, Eric Nyberg, and Vern R Walker. 2015. Introducing LUIIMA: an Experiment in Legal Conceptual Retrieval of Vaccine Injury Decisions Using a UIMA Type System and Tools. In *Proc. of the 15th Int. Conference on Artificial Intelligence and Law*. ACM, 69–78.
- [5] Guiraudé Lame. 2005. *Using NLP Techniques to Identify Legal Ontology Components: Concepts and Relations*. Springer Berlin Heidelberg, 169–184.
- [6] Alessandro Lenci, Simonetta Montemagni, Vito Pirrelli, and Giulia Venturi. 2007. NLP-based Ontology Learning from Legal Texts. A Case Study.. In *Proc. of the 2nd Workshop on Legal Ontologies and Artificial Intelligence Techniques*. Citeseer, 113–129.
- [7] José Saías and Paulo Quaresma. 2005. *A Methodology to Create Legal Ontologies in a Logic Programming Information Retrieval System*. Springer, 185–200.
- [8] Jaromir Savelka and Kevin D Ashley. 2016. Extracting Case Law Sentences for Argumentation about the Meaning of Statutory Terms. In *Proc. of the 3rd Int. Workshop on Argument Mining*. 50–59.
- [9] Jaromir Savelka, Matthias Grabmair, and Kevin D Ashley. 2014. Mining Information from Statutory Texts in Multi-Jurisdictional Settings. In *Proc. of the Int. Conference on Legal Knowledge and Information Systems*. IOS Press, 133–142.
- [10] The European Monitoring Centre for Drugs and Drugs Addiction. 2018. *Italy, Country Drug Report 2018*. Technical Report. The European Monitoring Centre for Drugs and Drugs Addiction.
- [11] Daniela Tiscornia. 2006. The LOIS project: Lexical Ontologies for Legal Information Sharing. In *Proc. of the V Legislative XML Workshop*. 189–204.
- [12] PRS Visser and TJM Bench-Capon. 1996. The Formal Specification of a Legal Ontology. In *Proc. of the Int. Conference on Legal Knowledge and Information Systems*.