

The Impact of Using Machine Learning for the Thematic Classification on Legal Documents

Aris Kosmopoulos

A.I. Researcher
SciFY PNPC and NCSR Demokritos
Athens, Greece
akosmo@scify.org

Stavroula Fikari

Attorney-at-law
Legal Informatics Consultant
Nomiki Bibliothiki
Athens, Greece
stavroula@nb.org

George Giannakopoulos

A.I. Researcher
SciFY PNPC and NCSR Demokritos
Athens, Greece
ggianna@iit.demokritos.gr

ABSTRACT

Gradually, the adaptation of Artificial Intelligence (AI) in various domains is becoming a fact. Although the legal domain offers several such opportunities, the ethical dilemmas that arise must be taken into serious consideration. In this work we demonstrate a real case scenario where the infusion of AI into a preexisting procedure can empower the human and facilitate the whole process of legal document annotation, as a supporting workflow related to legal AI. Furthermore, we discuss the ethical aspects of AI adoption, pointing out that the related ethical impact between different scenarios can vary greatly, offering the presented use case as an example of an AI application in the borderline of legal AI.

CCS CONCEPTS

• **Applied computing** → Law; **Annotation**; • **Computing methodologies** → **Supervised learning by classification**.

KEYWORDS

legal AI, document classification, multi-label classification, annotation

1 INTRODUCTION

An important aspect of Artificial Intelligence (AI) is the development of software that behaves and works like humans do. On the other hand, AI does not always try to replace humans. The facilitation of a human task is such a case, where AI can speed up the completion of an undertaken task and increase productivity.

AI can be applied in various domains and each domain naturally has certain characteristics and limitations that must be taken into account. Legal AI can refer to many different things, which can be grouped in two main categories:

- Legal issues arising from the use of AI systems similar to those arising from other innovative products and solutions and concerning the statutory and regulatory framework (data protection, consumers' rights, IP rights, competition).
- Employment of AI techniques and methods to produce tools and solutions assisting the legal professionals in every-day practice.

Although Legal AI offers several opportunities of AI applications, several ethical dilemmas must also be taken into consideration. For example allowing a computer program to create human laws, or even act as a judge, are indeed some very sensitive scenarios. But is this always case?

Facilitating the work of a human expert is a much less restrictive scenario in terms of ethical dilemmas. In this paper we focus on presenting a real-world application of AI in a legal setting. Nomiki Bibliothiki¹, a major legal content provider, has developed a website (legal content platform²) providing to legal professionals easy access to a full range of legal documents (legislation, case-law and other official legal documents, legal doctrine, templates of legal acts), which can support legal decision-making. A main concern is how the platform can arrange and classify this content in order to deliver quick, accurate and valid search results.

Among legal documents to be processed and analyzed are the administrative acts published in Issue B of the Official Government Gazette of Greece. A legal annotator must assign one or more subject-matter categories and legal terms chosen out of a hierarchical index (which is part of a thesaurus). The solution offered by AI – designed and implemented by SciFY PNPC³, an AI technology transfer and digital transformation not-for-profit company – was an automated classification process that proposed such categories and legal terms to the legal annotator. The benefit of this automated process is impressive and allows the annotator to perform the task much faster.

The contributions of this work are the following:

- The outline of a real-world use of AI use for legal domain tasks.
- A discussion on the benefits and the presence of ethical risks in this use case, but also a widening of the discussion to imply future, related concerns.

The rest of the document is structured as follows. In Section 2 we present some related work, while in Section 3 we describe the use case in more detail. In Section 4 we discuss some ethical aspects of the task and we conclude the paper in Section 5.

2 RELATED WORK

One of the essential steps in the analysis of large document collections is the thematic classification of these documents. As the volume of data increases significantly, manual analysis requires

¹<https://www.nb.org/>

²<https://www.qualex.gr>

³<https://www.scify.gr/>

effort and time. For that reason, over the last decades one of the main concerns of data scientists consists in designing processes of document analysis which tackle this challenge. Thus, they started experimenting with the implementation of automatic methods of classification. In this section, we refer to the most related applications of classification, since the literature of text classification in general is immense (cf. [1, 2, 7–9]).

In [4] a semi-automatic method, based on keyword classification of documents, assigns appropriate branches of knowledge to documents of Polish digital Libraries by using clustering algorithms. The experiment was conducted with the assistance of human annotators. The method was evaluated to be applicable to the thematic classification of documents in large digital collections.

An experiment of using machine learning (ML) techniques to classify sentences in Dutch legislation was used in [5]. These results are compared to the results of a pattern-based classifier and the conclusion was that pattern-based approach is preferable.

A domain specific approach regarding the classification of laws is presented in [3]. The system can compute similarities between small snippets of large heterogeneous laws. Another approach of classification and labeling of European laws is described in [11]. The authors state that the segmentation of each legal document into several parts can greatly improve the quality of labeling.

Another important consideration regarding legal document classification is whether linguistic information can help the classifiers. In [6] the authors evaluate the usefulness of adding lemmatization and part-of-speech in the classification pipeline and conclude that the results were in fact improved.

The limitations and perspectives of AI application in predictive justice was studied in [10]. The paper focuses on the Federal Court of Canada and examines the use of various state of the art methods of natural language processing and machine learning algorithms. Another case of application of AI in the legal domain is that of automatic summarization of legal texts. Such a goal was that of the SALOMON project presented in [12] that was applied to Belgian criminal cases.

In this work, we do not focus on the classification itself, but rather on the use case of classifying legal documents, as a support tool to efficient and effective legal content delivery. We also discuss ethical implications, but also the value added through the use of AI in this setting.

3 USE CASE DESCRIPTION

The so called “information crisis” in legal domain is a general phenomenon, meaning that the legal professionals need to access large volumes of legal information in order to treat a case and solve a legal problem. This crisis is aggravated by the diversity of legal sources to be consulted and, thus, the challenge mostly consists of locating and evaluating information delivered by various sources so as select and cite pertinent documents.

The content – provided to professionals by Nomiki Bibliothiki – to support legal decision-making, must always be indexed and classified in order to be delivered quickly and accurately. To accomplish this goal, several techniques of multi-level legal analysis are applied, like indexing and classification. But the rapidly increasing volume and complexity of data requires effort and time.

By 2018, the classification of the administrative acts published in Issue B of the Greek Official Government Gazette was being performed manually. The legal annotator was searching and choosing the relevant legal terms in a dedicated software tool (Figure 1), repeatedly for each term and for each separate legal act in two steps:

- Assignment of one or more subject-matter categories chosen out of a drop-down list.
- Assignment of legal terms chosen out of a hierarchical index (which is part of a thesaurus).

An AI solution of multi-label classification was designed and implemented by SciFY (Science For You). SciFY is a not-for-profit organization that implements digital transformation initiatives in the fields of Artificial Intelligence, assistive technologies, entrepreneurship, e-participation and education.

As in every machine learning training process the quality and quantity of data greatly affects the expected performance. This process was greatly facilitated by the excellent-quality, annotated data provided by Nomiki Bibliothiki. The provided legal document data were well structured and consistent, qualities ascertained by appropriate quality assurance processes. Another important factor for the success of the use case was the quantity of training instances per class, which in most cases was sufficient (in the order of tens of instances) in order to train a classification model.

For each class (categories and legal terms), given that sufficient training instances existed, we trained a binary classifier (approximately 1700 classifiers were used, one for each category / term). A bag-of-words approach was used in order to extract features from the legal documents (around 85 thousands of documents were used in total as training instances). A feature selection process was also applied to remove rare features and speed up the training and prediction processes, without negatively affecting the performance. During prediction, each instance (legal act) is evaluated by each classifier. When the classifier predicts with sufficient confidence that the document should be assigned the category label, the label is suggested to the human annotator as a plausible option (cf. Figure 2).

The performance of the suggestion is impressive: the internal tests on the actual workflow of the annotators showed a success rate of 98% (perceived estimated accuracy of the end user) in legal acts of standard and repetitive regulations. As a result, all the annotator has to do now, is accept all or part of the proposed terms in one move, instead of searching the drop-down lists.

We should notice, though, that the semi-automatic process of classification still remains a human-supervised method in order to avoid implied annotation risks (i.e. not using scarce classes which are not proposed by the algorithm) and the instruction given to annotators is to consider the addition of not proposed terms that are assessed as relevant or even to reject non pertinent proposed terms.

In any case, the time saved is significant, since for the standardized legal acts, which is the majority (almost 70%), the time annotation time was reduced by 50%. Given that, the legal annotators can now focus on more complex tasks of legal analysis, such as the consolidation of legal texts and the creation of links between related texts.

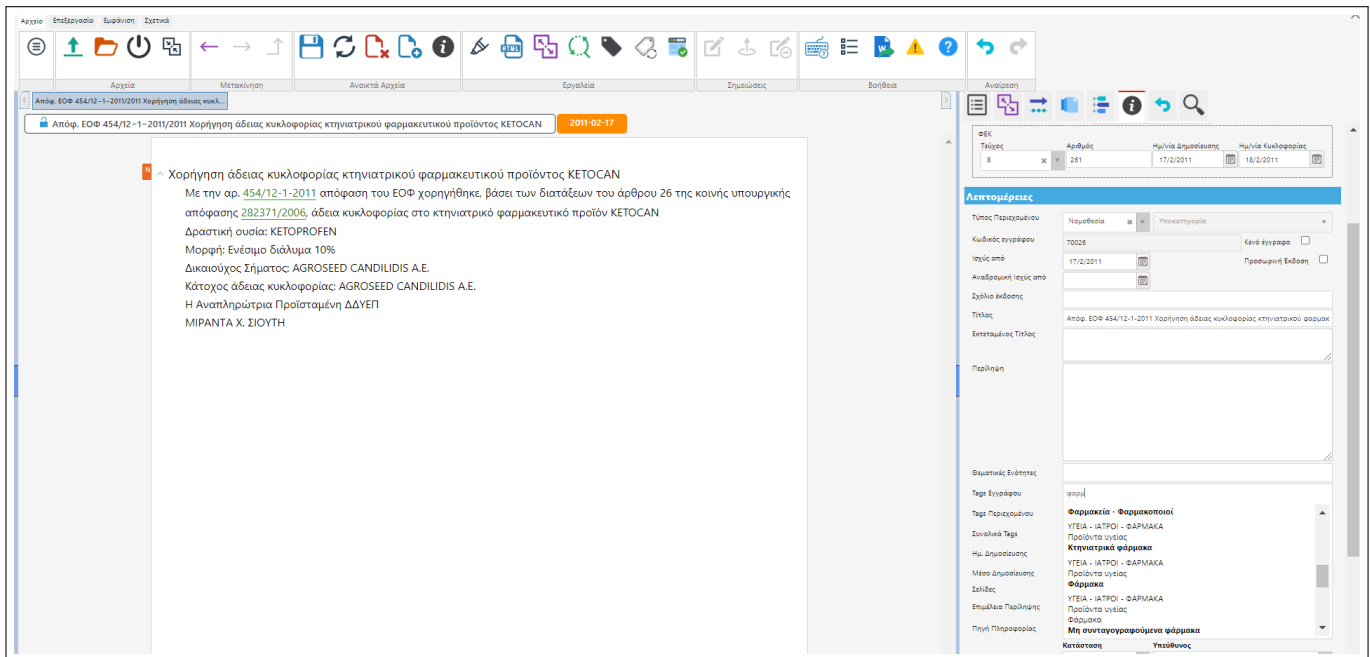


Figure 1: Tool used by a legal annotator in order to assign subject-matter categories and legal terms to a legal document.

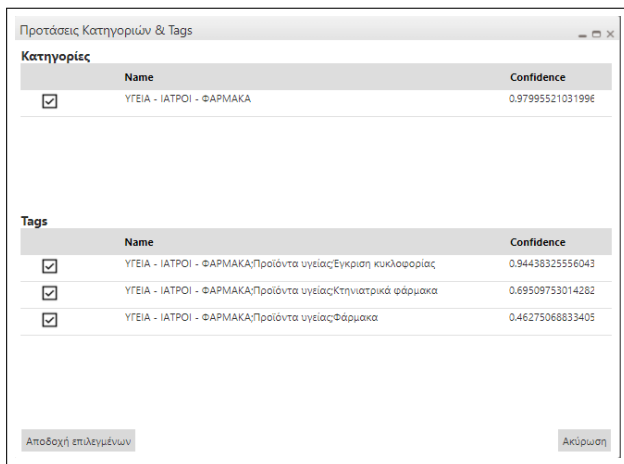


Figure 2: Legal annotator is facilitated by categories and terms predictions provided by AI.

Based on the above, the gain from the integration of AI components in the workflows of Nomiki Vivliothiki is clear. In the next section, we discuss ethical aspects of the system under the prism of legal AI ethical risks.

4 DISCUSSION OF THE ETHICAL ASPECTS

In the legal setting, there exist a number of subtle dangers in using AI, most notably:

algorithmic bias, which describes the preference that an algorithm may contain towards a specific decision. This bias

can be caused by inherent idiosyncrasies of an algorithm. This bias can be problematic in cases where the output of an algorithm implies or explicitly leads to a specific judicial outcome, e.g. a verdict.

data bias, which describes the implicit bias added to a machine learning algorithm, through the selection of training data. There exist several cases of such bias, again leading to unjust outcomes for a given legal setting.

explainability, which describes the danger of not being able to explain a decision of a machine learning system, while the decision significantly impacts a human subject. The usual reason for this risk related to the mathematical modeling of a problem in an AI system, which cannot provide a humanly-understandable response of the "why?" a decision was taken. The "explanation" is essentially a complex mathematical function, which may be impossible to interpret in meaningful terms.

default decisions, which refers to the danger of taking judicial decisions, without offering the possibility of rebuttal to the impacted subject.

agency and accountability of a decision, which refers to the challenge of assigning accountability to a person for a given decision, in the case when the decision was made by an AI system.

All the above challenges arise in cases where the legal process is directly affected by an AI supporting system. In this paper, however, we claim that there exist borderline applications of AI in the legal setting, where the above risks are mitigated. Essentially, these borderline applications refer to functions of AI in the information

gathering process, where there is always a human in the loop, and there exist at least two levels of validation for the AI outcomes.

In our use case, the AI system works to help the *annotation* of content *related to* legal settings. In other words, the AI is meant to help humans in increasing the indexability and retrievability of documents related to a legal setting. The AI decision is, thus, a suggestion to be validated by a human (the annotator) in the related quality assurance (QA) process. The results of this process allow legal professionals - the end users - to retrieve information related to their work, e.g. laws and decisions referring to similar cases. At this level, again a human is to finally decide what is related and what is not. Thus, the AI decisions are validated twice.

A hidden risk in this process is the fact that, once the end users increase their confidence towards the system, they may rely more and more on the document that the system retrieves. We consider the worst case scenario, where a critically related document was mistakenly classified by the system and, thus, is not retrieved as relevant to the end user query. It is possible that the outcome of the legal process is, thus, affected by the lack of this documentation.

Such a risk can be mitigated by two simple actions. The first relates to the validation of suggested classification tags by more than one human, minimizing the risk of erroneous tags. The second relates to the training of the end users, so that they utilize a minimum number of *different queries* to retrieve documents related to their case.

Cross-referencing the above discussion with the main identified risks of legal AI, we can see that:

- algorithmic and data bias is reduced through the quality assurance processes. Furthermore, even if there is bias, it does not directly affect the judicial processes, even though it may alter the flow of information towards the interested parties. In any case, the final decision still relies on humans.
- explainability may not be of real value in this setting, since the classification decision has limited impact and is easy to change, if the human annotator has a different opinion.
- the use of AI in our setting is not a part of the judicial processes themselves, but a supporting workflow for the gathering of related information.
- the agency and accountability of any decisions remains tied to the end user, who has always been responsible for the search and verification of gathered information.

Based on the above analysis, we suggest that such ethical/impact checklists could be useful to identify whether a given use case is a support process, as above, what are the related risks and how these risks can be mitigated.

In the following paragraphs, we go beyond the current use case we described, highlighting possible future directions of legal technology in Greece.

One possible future direction is that of Legal Research Solutions. Legal content providers use AI techniques to optimize legal research and deliver accurate results. The main features of such solutions are:

- The support of natural language search.
- The recognition of legal terminology.

- The analysis of legal documents through powerful citators, which allow the history tracking of a legal text and its treatment by official factors.
- The automatic summarization of documents.
- The production of litigation analytics.

Another direction is that of Predictive Analytics Solutions. AI tools utilize case law, public records, dockets, and jury verdicts to identify patterns in past and current data and then analyze the facts of a lawyer's case to provide an intelligent prediction of the outcome. Those tools can be extremely useful to legal practitioners and they are widely used in the USA and Canada. On the contrary in Europe there is a reticence due to ethical issues.⁴

Predictive Analytics methods can be applied to develop more advanced tools for legal risk assessment and legal risk management.

Legal risk can be defined in general as the risk of loss incurred to an organization or an individual due to factors related to legal issues. The various aspects of legal risk can be classified into the following broad categories:

- Litigation risk: potential legal disputes arising from business activities.
- Contractual risk: failure to fulfill contractual obligations by a contractual party resulting in liabilities and damages.
- Regulatory risk: modifications in legislation imposing new compliance practices and costs.
- Compliance risk: failure to comply with laws and regulations resulting in sanctions and penalties.

The legal uncertainty in the aspects mentioned above can affect a business or a market significantly and cause serious financial or other losses. The solutions and products offered use AI techniques which take into consideration and analyze legal data relevant to the circumstances of the person or entity concerned and assist them in developing an effective risk management strategy.

5 CONCLUSION

In this work we described a real-world application of AI in a legal setting. We showed how the infusion of AI into a pre-existing legal content generation process empowers the human and the requirements for such an application. We highlighted aspects of this empowerment in the use case and showed how a human-in-the-loop AI system can provide multiplicative effects to everyday work. We also described ethical aspects and challenges of the setting, but also of future prospects.

REFERENCES

- [1] Charu C Aggarwal and ChengXiang Zhai. 2012. A survey of text classification algorithms. In *Mining text data*. Springer, 163–222.
- [2] Berna Altunel and Murat Can Ganiz. 2018. Semantic text classification: A survey of past and recent advances. *Information Processing & Management* 54, 6 (2018), 1129–1153.
- [3] Guido Boella, Luigi Di Caro, and Llio Humphreys. 2011. Using classification to support legal knowledge engineers in the eunomos legal document management system. In *Fifth international workshop on Juris-informatics (JURISIN)*.
- [4] Łukasz Borchmann, Filip Gralinski, Rafał Jaworski, and Piotr Wierzechon. 2015. A semi-automatic method for thematic classification of documents in a large text corpus. *Corpus-Based Research in the Humanities (CRH)* (2015), 13.
- [5] Emile de Maat, Kai Krabben, Radboud Winkels, et al. 2010. Machine Learning versus Knowledge Based Classification of Legal Texts.. In *JURIX*. 87–96.

⁴See the case of France: statutory prohibition of court decisions analysis based on the judge profile.

- [6] Teresa Gonçalves and Paulo Quaresma. 2005. Is linguistic information relevant for the classification of legal texts?. In *Proceedings of the 10th international conference on Artificial intelligence and law*. 168–176.
- [7] Vandana Korde and C Namrata Mahender. 2012. Text classification and classifiers: A survey. *International Journal of Artificial Intelligence & Applications* 3, 2 (2012), 85.
- [8] Kamran Kowsari, Kiana Jafari Meimandi, Mojtaba Heidarysafa, Sanjana Mendu, Laura Barnes, and Donald Brown. 2019. Text classification algorithms: A survey. *Information* 10, 4 (2019), 150.
- [9] Nikiforos Pittaras, George Giannakopoulos, George Papadakis, and Vangelis Karkaletsis. 2020. Text classification with semantically enriched word embeddings. *Natural Language Engineering* (2020), 1–35.
- [10] Marc Queudot and Marie-Jean Meurs. 2018. Artificial intelligence and predictive justice: Limitations and perspectives. In *International Conference on Industrial, Engineering and Other Applications of Applied Intelligent Systems*. Springer, 889–897.
- [11] Erich Schweighofer, Andreas Rauber, and Michael Dittenbach. 2001. Automatic text representation, classification and labeling in European law. In *Proceedings of the 8th international conference on Artificial intelligence and law*. 78–87.
- [12] Caroline Uyttendaele, Marie-Francine Moens, and Jos Dumortier. 1998. Salomon: automatic abstracting of legal cases for effective access to court decisions. *Artificial Intelligence and Law* 6, 1 (1998), 59–79.