# Data Science Applications and Implications in Legal Studies: A Perspective Through Topic Modelling

Jinzhe Tan[1], Huan Wan[1], Ping Yan[1], and Zhen Zhu[2,*]

[1]*Zhongnan University of Economics and Law, Wuhan, China*
[2]*University of Kent, Canterbury, United Kingdom*

## Abstract

Law and legal studies has been an exciting new field for data science applications whereas the technological advancement also has profound implications for legal practice. For example, the legal industry has accumulated a rich body of high quality texts, images and other digitised formats, which are ready to be further processed and analysed by data scientists. On the other hand, the increasing popularity of data science has been a genuine challenge to legal practitioners, regulators and even general public and has motivated a long-lasting debate in the academia focusing on issues such as privacy protection and algorithmic discrimination. This paper collects 1236 journal articles involving both law and data science from the platform Web of Science to understand the patterns and trends of this interdisciplinary research field in terms of English journal publications. We find a clear trend of increasing publication volume over time and a strong presence of high-impact law and political science journals. We then use the Latent Dirichlet Allocation (LDA) as a topic modelling method to classify the abstracts into four topics based on the coherence measure. The four topics identified confirm that both challenges and opportunities have been investigated in this interdisciplinary field and help offer directions for future research.

**Keywords** *artificial intelligence; law; literature review; text mining*

## 1 Introduction

With the rapid development of the Internet and information technologies, varieties of data are being generated at an unprecedented scale across industries. Massive data resources have gradually become the basic resources of countries and societies, pushing the world into the era of Big Data. In response to this, the legal industry has started to vigorously promote the electronic, digitised and publicised legal information, making Legal Big Data a new buzzword. The emergence of Legal Big Data helps data science find applications in law and brings revolutionary changes to legal research methods and practice (Frankenreiter and Livermore, 2020). At present, technologies in the fields of artificial intelligence, natural language processing, text mining, network analysis and machine learning are increasingly adopted by legal practitioners and legal scholars (Aletras et al., 2016; Dhami and Belton, 2016; Medvedeva et al., 2020; Olsen and Küçüksu, 2017; Tarissan and Nollez-Goldbach, 2016; Wyner et al., 2010). These new tools are used to solve not only practical challenges in judicial practice and legal disciplines, but also academic research problems related to law and legal systems.

*Corresponding author, Kent Business School, University of Kent, Parkwood Road, Canterbury, Kent CT2 7FS, United Kingdom. Tel: +44 (0)1227 827726. Email: z.zhu@kent.ac.uk.

Despite the many benefits to be realised, data science also has the potential to significantly reshape the legal industry and challenge the existing legal systems. We therefore divide the literature involving both data science and law into two topics below: one is on the applications of data science in law and the other is on the implications of data science for legal practice.

## 1.1 Applications of Data Science in Law

The applications of data science can be roughly summarised in two fields (Frankenreiter and Livermore, 2020). One is to study law as code, for example, the conversion of legal rules into computer code. In recent decades, legal scholars and computer scientists have tried to develop computer codes that represent legal rules in regulations and case law. This research method regards law as a set of formal rules, which can eventually be formalised as inputs that can be directly processed by a computer. Some even optimistically discuss that in the future, legislators may transform common natural language in laws and regulations into executable codes in public programming languages. This language would be clearer and could be customised more easily according to specific situations (Fagan and Levmore, 2019).

The other field is to study law as data, which is based on the rich data resources generated in the legal processes. The ever-expanding legal system is nowadays a digital environment with a comprehensive record of traces, which is highly suitable for researchers to carry out Big Data analysis and experiments. This field deserves a more detailed review in this paper as it is closely related to data science. For example, legal texts can be digitally translated into data and quantitatively analysed. The use of data analysis methods complements the traditional normative legal research and provides a new direction for exploring legal issues based on facts and data (Frankenreiter and Livermore, 2020).

As a concrete example, in the United States, there is a tradition of quantitative analysis of case law. There are some quantitative studies on the case law data sets of US courts. Most of these studies use manually collected and coded case law. Many studies use the Supreme Court database, which contains manually collected and professionally coded data on the conduct of the US Supreme Court over the past 200 years (Katz et al., 2017). This type of research analyses the relationship between the gender or political background of judges and their decision-making and similar studies have been carried out for the International Criminal Court (Tarissan and Nollez-Goldbach, 2016), the European Court of Justice (Frankenreiter, 2017) and the European Court of Human Rights (Madsen, 2018).

In terms of data science tools, many studies provide basic descriptive statistics of manually collected and coded case law (Madsen, 2018). Some studies have shown relatively basic correlation analysis or statistical test results (Bruijn et al., 2018) whereas others have employed more complex statistical analyses, including regression analysis of case law (Dhami and Belton, 2016). Noticeably a growing number of studies have applied network analysis to legal data (Olsen and Küçüksu, 2017; Šadl and Olsen, 2017; Tarissan and Nollez-Goldbach, 2016). Text mining and machine learning have also gained popularity in recent studies (Aletras et al., 2016; Katz et al., 2017; Medvedeva et al., 2020; Wyner et al., 2010).

In short, in the era of Big Data, the electronic records of law and its implementations have accumulated a large amount of data. The increase in data types and volume, the improvement in data quality and the advancement of data processing capabilities have created an excellent environment for quantitative legal analysis. Data science, by providing computational methods, has facilitated important contributions in a diverse set of law-related research areas. As these tools continue to advance, and law scholars become more familiar with their potential applications,

the impact of data science on law is likely to continue to grow (Frankenreiter and Livermore, 2020).

## 1.2 Implications of Data Science for Legal Practice

Data science has undoubtedly facilitated the development of legal research and practice with its efficient processes. However, the gradual adaptation and acceptance of data science in legal practice has been met with many challenges (Wendel, 2019). The impact of data science on law comes from two directions: externally, the applications of data science in other fields have an impact on the practice of law and requires a legal response; internally, the applications of data science in the legal field itself also have an impact on law and require corresponding changes in law.

### 1.2.1 The External Impact of Data Science on Legal Practice

Ever since the Internet era, data sources have become more and more extensive. For example, purchase records, app usage, and other personal data are used to form user profiles for commercial purposes (Tene and Polonetsky, 2011). Banks and other financial institutions have access to customers' intention to borrow money and their repayment history to decide whether to grant them a loan or not (Martin, 2019). The rules of these scoring systems are so complex, and the parameters considered are so numerous, that customers may be completely unaware of which of their actions affect their credit scores. As the scale of data collection continues to grow, some scoring systems have evolved into untouchable 'tyrants' where points are deducted for simply trying to figure out the scoring mechanism of the scoring system (Citron and Pasquale, 2014). These improper collections and uses of personal information require timely legal responses. The enactment of laws such as the EU's GDPR and China's Personal Information Protection Law represents the legal community's response to the challenges faced in the era of Big Data.

With the advancement of information technologies such as deep forgery technologies and neural networks, artificial intelligence forgery technologies represented by "deep fake" has become increasingly sophisticated. The "deep fake" technology can use a portion of the original material to forge a person's handwriting, voice, or face replacement in a video to achieve the effect of faking. This technology has applications in education, arts, and self-governance, but the risks associated with its misuse are more significant, with potential threats to personal privacy, media trust, economic development, and public safety, posing significant challenges to legal practice (Chesney and Citron, 2019).

### 1.2.2 The Internal Impact of Data Science on Legal Practice

The importance of data science is also growing rapidly within the legal practice systems and academic. Scholars are exploring the integration of data science with the law itself, and courts are experimenting with the integration of data science into judicial trials, such as China's smart courts (Zheng, 2020), which expect to use data science techniques to free lawyers and judges from the drudgery of their work. Universities and academic institutions are also experimenting with the integration of data science and law to train people who can adapt to future directions (Hod et al., 2022).

Some data science applications closely related to law practice have been developed in recent years. For example, Quemy and Wrembel (2020) used the European Court of Human Rights (ECtHR) database for the outcome identification task, achieved the average accuracy of 94.43%.

For the task of automated court decision prediction, Medvedeva et al. (2021) conducted an analysis of the ECtHR database, which resulted in a maximum prediction accuracy of 66%. However, many scholars still point out that data science can only play an auxiliary role for judges and lawyers to analyse cases, because algorithms only analyse structured data according to complex procedures, and cannot understand the actual meaning of cases and judgments. Machine learning gradually converges to human analysis with continuous training and also represents an increase in the complexity of the internal system to the extent that the observer can only know the input and output sides of the algorithm and not the internal computational process of the algorithm or whether there is potential discrimination in the computational process (Hildebrandt, 2018).

The challenges of deploying data science applications in law practice are often associated with potential discrimination and the low interpretability and transparency of the models. Existing legal data are not completely non-discriminant, and potential biases by legal practitioners are implicit in legal data sets (Tene and Polonetsky, 2017; Surden, 2019). On the other hand, the low interpretability and transparency of the data science models are due to the increasing complexity of the internal system, where an observer can only know the input and output side of the algorithm, but not the internal computational process of it (Bibal et al., 2021).

These criticisms show that this interdisciplinary research area is still in its infancy. Sunstein (2001) believes that the current integration of data science with law is only an advanced version of LEXIS. There is only a quantitative change, not a qualitative one, when compared with the traditional legal research methods. Moreover, Wendel (2019) argues that even in a highly technologically advanced future, data science models can only simulate human moral reasoning but can never be a free and equal person.

The summary above is likely to be biased as it is based on our domain knowledge and preferences. Therefore, we also take a systematic approach based on text mining to generate an exploratory literature review. The rest of the paper is organised as follows: Section 2 describes our literature search strategies and the text mining methods used for our exploratory literature review, Section 3 presents the results, and finally Section 4 concludes the paper.

## 2  Methodology

We use the platforms Web of Science (webofscience.com) for literature search and data collection. We restrict the categories to be "law" and "political science", the document type to be "article", and the language to be "English". To retrieve the literature involving both law and data science, our search statement requires that at least one out of the terms "law" and "legal" (i.e., related to law) *and* at least one out of the terms "data science", "machine learning", "artificial intelligence", "big data", "data mining", "text mining", "natural language processing", "quantitative", and "computational" (i.e., related to data science) should be present in article title, abstract or keywords. Furthermore, the articles need to be indexed by SSCI (Social Sciences Citation Index), ESCI (Emerging Sources Citation Index), or SCI-EXPANDED (Science Citation Index Expanded).

Traditionally being manual and tedious, exploratory literature reviews are nowadays automated with the advancement of text mining methods (Asmussen and Møller, 2019). Among others, topic modelling has been an increasingly popular method of summarising literature by classifying a large body of documents into topics (Asmussen and Møller, 2019; Li and Lei, 2021). The default method of topic modelling is the Latent Dirichlet Allocation (LDA) (Blei et al., 2003), which is an unsupervised learning method and therefore does not require a high

upfront cost of time or domain knowledge (Asmussen and Møller, 2019). The output of LDA algorithms includes two distributions, the topic-word distribution $\mathbf{t}_i = \{t_{ij}\}$, where $i$ denotes the topic and $j$ denotes the word and $t_{ij}$ denotes the relative weight of word $j$ in topic $i$, and the document-topic distribution $\mathbf{d}_s = \{d_{si}\}$, where $i$ again denotes the topic and $s$ denotes the document and $d_{si}$ denotes the relative weight of topic $i$ in document $s$.

To our awareness, the most closely related works to our paper are by Wieringa (2020) and by Rosca et al. (2020). Whereas both of them are exploratory literature reviews in related fields from a text mining approach, our paper differs from them in important ways. Wieringa (2020) focused on the topic of "algorithmic accountability" and did a network visualisation of the keywords. Our paper, however, considers the much broader domains of "data science" and "law" and uses the topic modelling method to summarise the literature. Rosca et al. (2020) focused on the topic of "artificial intelligence" and collected the articles from the HeinOnline database specialised in legal studies. Our paper, on the other hand, retrieves our sample from the cross-discipline platform Web of Science and compares two mainstream LDA algorithms, gensim (Řehůřek and Sojka, 2010) and mallet (McCallum, 2002), in order to determine the number of topics. The difference between gensim and mallet is that gensim uses the variational Bayes sampling method which is faster but less precise than the Gibbs sampling method used by mallet.

One key parameter in LDA algorithms is the number of topics, often denoted as $k$. We use the coherence measure to determine $k$. In simple terms, topic coherence measures whether words in a topic tend to co-occur together (Röder et al., 2015).

## 3   Results

After implementing the search strategies above and removing the items missing abstracts or publication years, we have collected 1236 articles from the Web of Science. As shown in Figure 1, the first article meeting our search criteria was published in the year as early as 1993 (Miller, 1993). It was then followed by almost two decades of inactivity and a dramatic growth since 2011.

There are 360 distinct journals included in our sample. We identify the high-impact journals by taking into account the number of citations. It turns out that there are 96 journals with articles that are cited by at least 20 times in our sample. Figure 2 summarise the distributions of the articles over these 96 journal titles. They include a diverse set of journals in the categories of law, political science, and interdisciplinary. The interdisciplinary research topic between data science and law has made its presence known in the most prestigious law journals such as *Journal of Empirical Legal Studies*, *Harvard Law Review*, *Yale Law Journal* and *Stanford Law Review*, as well as political science journals such as *Journal of Politics*, *American Political Science Review* and *American Journal of Political Science*. On the other hand, Figure 2 are clearly dominated by two dedicated interdisciplinary journals, *Computer Law & Security Review* and *Artificial Intelligence and Law*.

We then use the LDA topic modelling method to detect the topics in our sample of abstracts. First, we compare the two mainstream LDA algorithms, gensim (Řehůřek and Sojka, 2010) and mallet (McCallum, 2002), and use the coherence measure to determine the number of topics $k$. For each number of topics from 2 to 15, we run each algorithm 10 times and then obtain the average coherence scores. The average coherence scores for both algorithms, along with the 95% confidence interval, are shown in Figure 3. Clearly mallet outperforms gensim across the different values of $k$. We choose $k = 4$ as it has the maximum average coherence score.
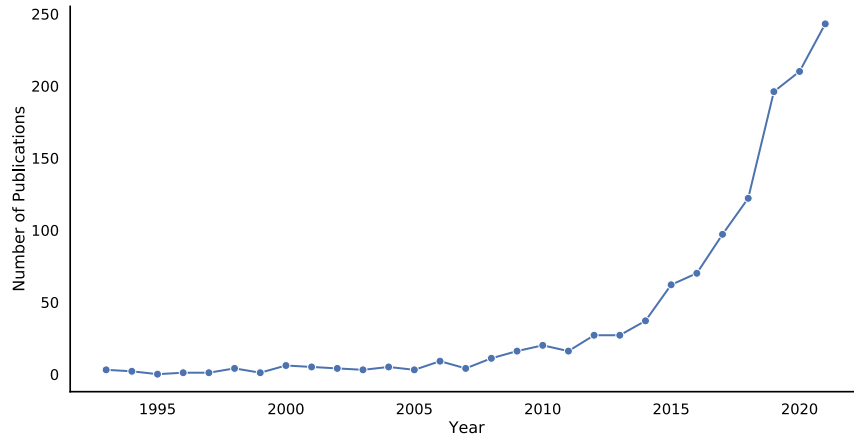
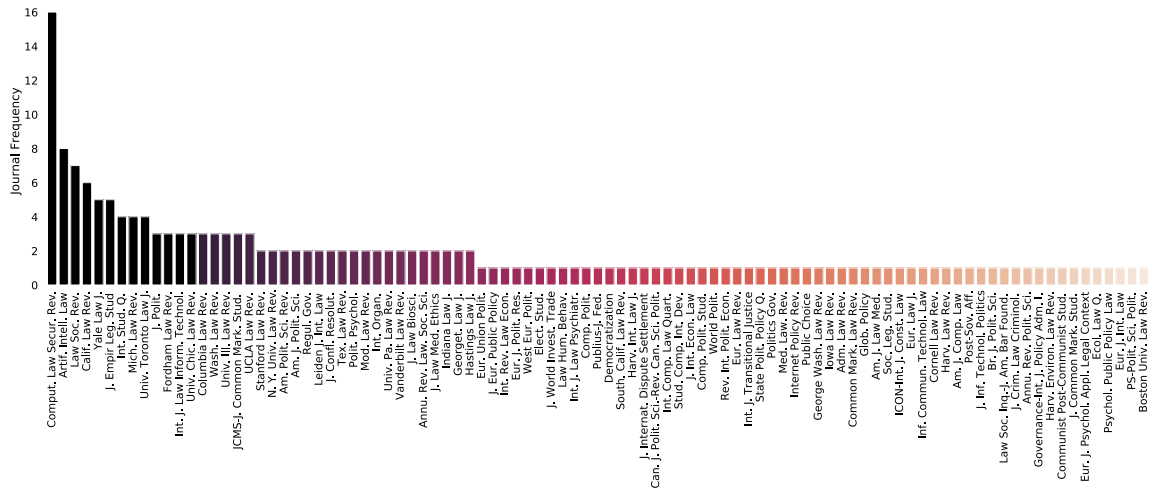Figure 1: Number of articles over time.


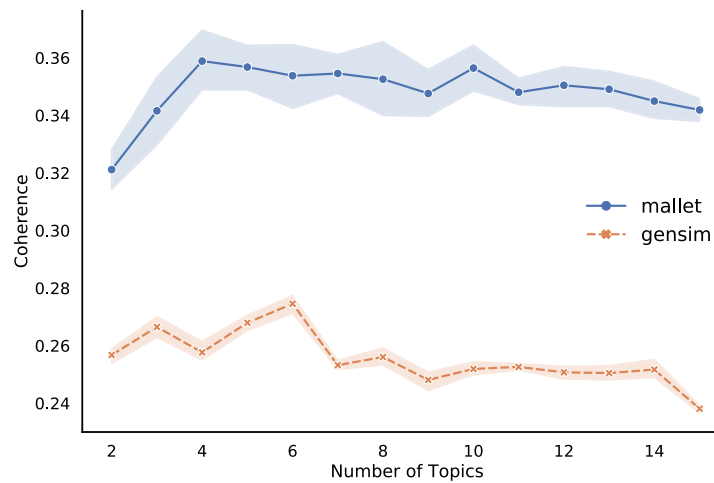
Figure 2: Frequency of high-impact journals.



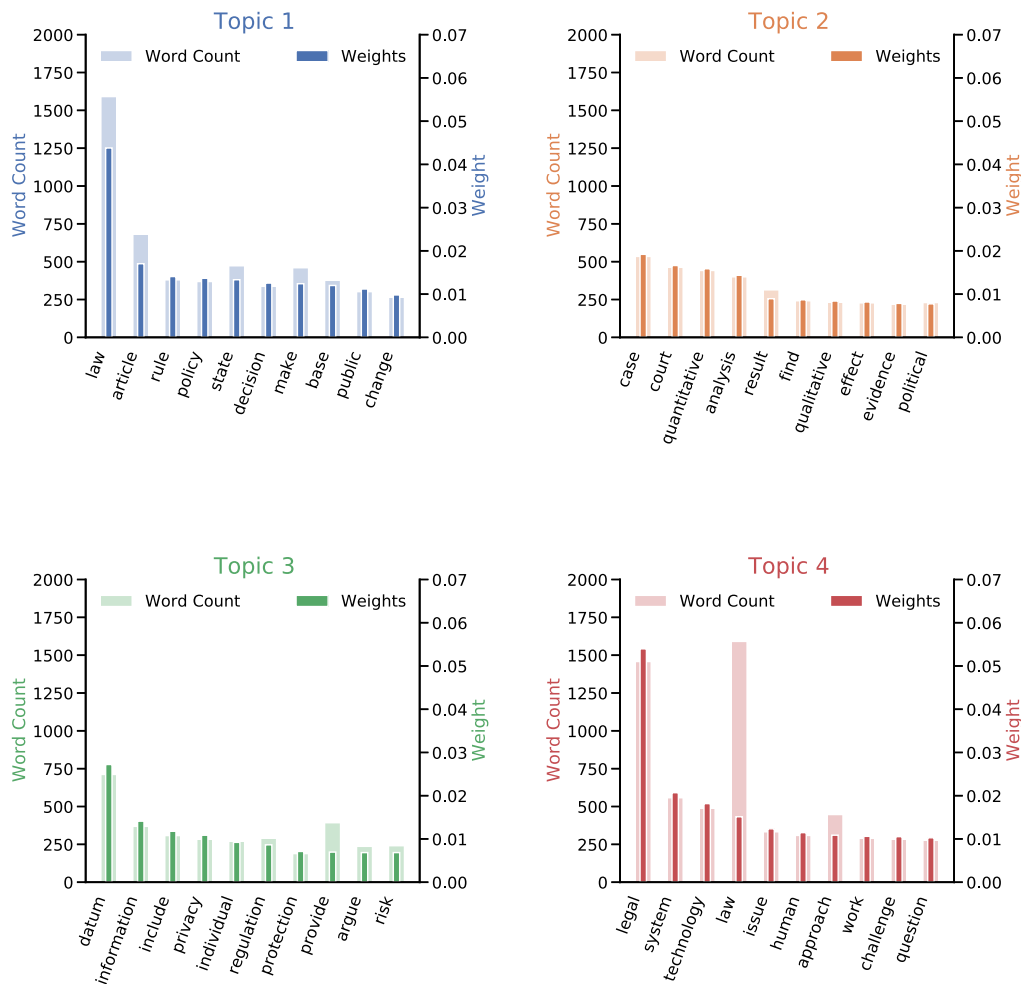Figure 3: Average coherence by number of topics.

Figure 4: Keywords.

We then visualise the results by showing the most weighted keywords and top frequent words ("wordcloud") for each topic automatically identified. Figure 4 and Figure 5 show the keywords and "wordcloud" for each topic respectively. We label the four topics as, Topic 1 "policy and governance" (also see the words such as "public", "government", and "state" in Figure 5), Topic 2 "methods and measures" (also see the words such as "quantitative", "empirical", and "analysis" in Figure 5), Topic 3 "data and privacy" (also see the words such as "datum", "information", and "protection" in Figure 5), and Topic 4 "system and technology" (also see the words such as "ai", "machine", and "digital" in Figure 5). These four topics are obviously interrelated, but each of them is focused on a slightly different angle, as indicated by their labels above.

These topics can be aggregated to the two discussed in Section 1 with our domain knowledge. That is, Topic 1 "policy and governance", Topic 3 "data and privacy", and Topic 4 "system and technology" mostly belong to the implications of data science for legal practice (see Section 1.2), whereas Topic 2 "methods and measures" mostly belongs to the applications of data science in law (see Section 1.1). Moreover, within the implications, Topic 3 is clearly more related to the external impact of data science on legal practice (see Section 1.2.1), whereas Topic 1 and Topic 4 possibly also cover the internal impact of data science on legal practice (see Section 1.2.2).

Figure 5: "Wordcloud".



Figure 6: Dominant topics distribution.

Therefore, the automated result of topic modelling supports the previous manual summary of the literature based on our domain knowledge. Importantly, our result differs from those of previous studies in that we are able to uncover Topic 2 "methods and measures" on the applications of data science in law by using the cross-discipline platform Web of Science instead of specialised databases such as HeinOnline (Rosca et al., 2020).

To understand the current trend as well as future directions in this interdisciplinary field, we assign each article to its dominant topic (i.e., $argmax_i d_{si}$) and present the number of articles for each topic in Figure 6. It turns out that the most popular topic is Topic 2 labelled as "methods and measures", which implies that the methodological innovations inspired by data science are well received in the legal research community.

## 4   Conclusions

The relations between data science and law are characterised by opportunities as well as challenges, as evidenced by our exploratory literature review above. On the one hand, law and legal

studies has always been a data-rich field where data science tools can be readily applied. On the other hand, data science technologies have the potential to reshape the legal profession and have raised many challenging issues such as how privacy can be preserved and how algorithmic discrimination needs to be regulated.

This paper systematically examines this still nascent interdisciplinary field in terms of research publications. Using our domain knowledge, we first divide the literature roughly into two topics, applications of data science in law and implications of data science for legal practice. We then collect the journal articles involving both law and data science from the Web of Science. We only consider the articles indexed by SSCI, ESCI, or SCI-EXPANDED. As a result, we have collected 1236 articles and 360 unique journals. We then take a text mining approach to present an automated exploratory literature review. There is a clear trend of increasing publication volume over time and of an increasing presence in mainstream law and political science journals.

Moreover, we use the Latent Dirichlet Allocation (LDA) as a topic modelling method to automatically detect the topics in our sample. The four topics detected are respectively, Topic 1 labelled as "policy and governance", Topic 2 labelled as "methods and measures", Topic 3 labelled as "data and privacy", and Topic 4 labelled as "system and technology". Topic 1, Topic 3, and Topic 4 mostly belong to the implications of data science for legal practice (see Section 1.2), whereas Topic 2 mostly belongs to the applications of data science in law (see Section 1.1). Our contribution is threefold. First, we ensure the quality of our sample by using a set of keywords highly related to data science and law and by using the cross-discipline platform Web of Science rather than databases specialised in legal studies or data science only. Second, we use the LDA topic modelling method to automatically summarise the literature involving both data science and law and compare the result with the manual summary based on our domain knowledge. Third, we carefully decide the number of topics by comparing two mainstream algorithms of LDA based on the coherence measure. The four topics identified support our intuition and confirm that both challenges and opportunities have been investigated in this interdisciplinary field and help offer directions for future research.

Being an exploratory literature review, this paper has several limitations. First, the topics identified by the LDA could be biased by the fact that we have restricted our search in the categories of law and political science. This may explain why the most frequent words in Figure 5 seem to be more on the implications of data science for legal practice rather than on the applications of data science in law. However, allowing the categories to go beyond law and political science will likely result in noisy data if we only consider article title, abstract and keywords. To gain a better understanding of the applications of data science in law, one possible solution is to focus on the specialised journals such as *Artificial Intelligence and Law*. Second, currently we consider all the articles' abstracts equally when implementing the text mining methods. The robustness of the results could be further checked if the abstracts are weighted by, for instance, article citations or journal impact factors. These limitations and challenges will have to await future research.

## Supplementary Material

The file "JDS_dataScienceLaw.ipynb" has the Python code used for the analysis above. The file "articles_en.csv" has the original data collected from Web of Science. The file "README.txt" has the description of the two files above.

# References

Aletras N, Tsarapatsanis D, Preoţiuc-Pietro D, Lampos V (2016). Predicting judicial decisions of the European Court of Human Rights: A natural language processing perspective. *PeerJ Computer Science*, 2: e93.

Asmussen CB, Møller C (2019). Smart literature review: A practical topic modelling approach to exploratory literature review. *Journal of Big Data*, 6(1): 1–18.

Bibal A, Lognoul M, De Streel A, Frénay B (2021). Legal requirements on explainability in machine learning. *Artificial Intelligence and Law*, 29(2): 149–169.

Blei DM, Ng AY, Jordan MI (2003). Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3: 993–1022.

Bruijn LM, Vols M, Brouwer JG (2018). Home closure as a weapon in the Dutch war on drugs: Does judicial review function as a safety net? *International Journal of Drug Policy*, 51: 137–147.

Chesney B, Citron D (2019). Deep fakes: A looming challenge for privacy, democracy, and national security. *California Law Review*, 107: 1753.

Citron DK, Pasquale FA (2014). The scored society: Due process for automated predictions. *Washington Law Review*, 89: 1.

Dhami MK, Belton I (2016). Statistical analyses of court decisions: An example of multilevel models of sentencing. *Law and Method*, 10: 247–266.

Fagan F, Levmore S (2019). The impact of artificial intelligence on rules, standards, and judicial discretion. *Southern California Law Review*, 93: 1.

Frankenreiter J (2017). The politics of citations at the ECJ—policy preferences of EU member state governments and the citation behavior of judges at the European Court of Justice. *Journal of Empirical Legal Studies*, 14(4): 813–857.

Frankenreiter J, Livermore MA (2020). Computational methods in legal analysis. *Annual Review of Law and Social Science*, 16: 39–57.

Hildebrandt M (2018). Law as computation in the era of artificial legal intelligence: Speaking law to the power of statistics. *University of Toronto Law Journal*, 68(supplement 1): 12–35.

Hod S, Chagal-Feferkorn K, Elkin-Koren N, Gal A (2022). Data science meets law. *Communications of the ACM*, 65(2): 35–39.

Katz DM, Bommarito MJ, Blackman J (2017). A general approach for predicting the behavior of the Supreme Court of the United States. *PLOS ONE*, 12(4): e0174698.

Li X, Lei L (2021). A bibliometric analysis of topic modelling studies (2000–2017). *Journal of Information Science*, 47(2): 161–175.

Madsen MR (2018). Rebalancing European human rights: Has the Brighton Declaration engendered a new deal on human rights in Europe? *Journal of International Dispute Settlement*, 9(2): 199–222.

Martin K (2019). Ethical implications and accountability of algorithms. *Journal of Business Ethics*, 160(4): 835–850.

McCallum AK (2002). Mallet: A machine learning for language toolkit. http://mallet.cs.umass.edu.

Medvedeva M, Üstun A, Xu X, Vols M, Wieling M (2021). Automatic judgement forecasting for pending applications of the European Court of Human Rights. In: *Proceedings of the Fifth Workshop on Automated Semantic Analysis of Information in Legal Text (ASAIL 2021)*.

Medvedeva M, Vols M, Wieling M (2020). Using machine learning to predict decisions of the European Court of Human Rights. *Artificial Intelligence and Law*, 28(2): 237–266.

Miller AR (1993). Copyright protection for computer programs, databases, and computer-generated works: Is anything new since CONTU? *Harvard Law Review*, 106(5): 978–1073.

Olsen HP, Küçüksu A (2017). Finding hidden patterns in ECtHR's case law: On how citation network analysis can improve our knowledge of ECtHR's Article 14 practice. *International Journal of Discrimination and the Law*, 17(1): 4–22.

Quemy A, Wrembel R (2020). On integrating and classifying legal text documents. In: *International Conference on Database and Expert Systems Applications*, 385–399. Springer.

Řehůřek R, Sojka P (2010). Software framework for topic modelling with large corpora. In: *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, 45–50. ELRA, Valletta, Malta.

Röder M, Both A, Hinneburg A (2015). Exploring the space of topic coherence measures. In: *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining*, 399–408.

Rosca C, Covrig B, Goanta C, van Dijck G, Spanakis G (2020). Return of the AI: An analysis of legal research on artificial intelligence using topic modeling. In: *NLLP@ KDD*, 3–10.

Šadl U, Olsen HP (2017). Can quantitative methods complement doctrinal legal studies? using citation network and corpus linguistic analysis to understand international courts. *Leiden Journal of International Law*, 30(2): 327–349.

Sunstein CR (2001). Of artificial intelligence and legal reasoning. *The University of Chicago Law School Roundtable*, 8(1): 29–35.

Surden H (2019). Artificial intelligence and law: An overview. *Georgia State University Law Review*, 35: 19–22.

Tarissan F, Nollez-Goldbach R (2016). Analysing the first case of the International Criminal Court from a network-science perspective. *Journal of Complex Networks*, 4(4): 616–634.

Tene O, Polonetsky J (2011). Privacy in the age of big data: A time for big decisions. *Stanford Law Review*, 64: 63.

Tene O, Polonetsky J (2017). Taming the golem: Challenges of ethical algorithmic decision-making. *North Carolina Journal of Law and Technology*, 19: 125.

Wendel WB (2019). The promise and limitations of artificial intelligence in the practice of law. *Oklahoma Law Review*, 72: 21.

Wieringa M (2020). What to account for when accounting for algorithms: A systematic literature review on algorithmic accountability. In: *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, 1–18.

Wyner A, Mochales-Palau R, Moens MF, Milward D (2010). Approaches to text mining arguments from legal cases. In: *Semantic Processing of Legal Texts: Where the Language of Law Meets the Law of Language* (E Francesconi, S Montemagni, W Peters, D Tiscornia, eds.), 60–79. Springer, Berlin Heidelberg.

Zheng GG (2020). China's grand design of people's smart courts. *Asian Journal of Law and Society*, 7(3): 561–582.