# Untangling the GDPR Using ConRelMiner

Karolin Winter, Stefanie Rinderle-Ma

Faculty of Computer Science, University of Vienna, Vienna, Austria

{karolin.winter, stefanie.rinderle-ma}@univie.ac.at

**Abstract**

The General Data Protection Regulation (GDPR) poses enormous challenges on companies and organizations with respect to understanding, implementing, and maintaining the contained constraints. We report on how the ConRelMiner method can be used for untangling the GDPR. For this, the GDPR is filtered and grouped along the roles mentioned by the GDPR and the reduction of sentences to be read by analysts is shown. Moreover, the output of the ConRelMiner – a cluster graph with relations between the sentences – is displayed and interpreted. Overall the goal is to illustrate how the effort for implementing the GDPR can be reduced and a structured and meaningful representation of the relevant GDPR sentences can be found.

# 1 Introduction and the ConRelMiner Method

Providing support for analyzing regulatory documents is of utmost importance for many companies nowadays as they face constantly changing or new requirements such as recently the General Data Protection Regulation (GDPR)[1]. Nowadays this is mostly done in a manual way which can be error-prone and costly. Hence, our recent research (cf.[11, 12]) aims at providing (semi-)automatic means to analyze regulatory documents based on text and data mining methods.

We aim at facilitating the handling of complicated and extensive regulatory documents such as legal texts. Therefore, we have developed a method, that is able to structure the documents accordingly and to detect relations between sentences. As a result, similar sentences are highlighted which reduces the reading effort. In addition, a grouping based on, e.g., given topics serves to identify text passages that are relevant for a specific user. Figure 1 outlines the basic idea. In addition, already implemented documents can also be integrated in order to detect conflicts and (partly) overlaps between sentences stemming from recently added documents. This supports the analysis of evolving regulatory documents over time.

The ConRelMiner method [11] consists of three steps, i.e., pre-processing, processing, and post-processing.

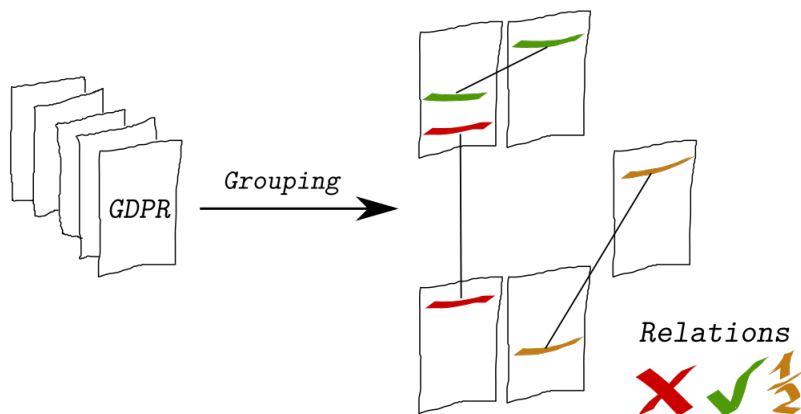---

[1] https://eugdpr.org/

Figure 1: Extraction and Grouping of Constraints from Regulatory Documents

The pre-processing step contains typical steps such as stemming and removal of stopwords. Novel is the fragmentation of the documents as described in [12] where documents can be split along a certain semantics, e.g., paragraphs. As shown in [12] this already enables a characterization of the documents, i.e., it can be derived which paragraph is associated with which theme or topic. Then the sentences that contain constraints are filtered out by using signal words.

The processing step employs techniques from text mining [2] and Natural Language Processing (NLP) [9], but also comprises novel concepts such as grouping sentences along topics and determining relations between the sentences based on their similarity. The grouping of sentences, resp. constraints can be customized individually depending on the type or size of the documents as well as additionally available information. In particular, a user can chose from three different methods. The first method uses term frequencies, whereas the second one exploits the structure of sentences and the third enables the integration of domain knowledge. Currently supported relations are "redundant", "subsumed", and "conflicting". For relations "redundant" and "subsumed" the related constraints can be viewed together and merged where applicable. Conflicting constraints can be also of interest, for example, if constraints contradict corresponding constraints in previous versions.

The result of the ConRelMiner method is a graph, which reflects the ordering by topics. In addition, sentences (the nodes of the graph) can be connected by edges describing the relations between them. Sentences can be "redundant" (marked green in the first figure, labeled with r in the graph), "subsumed" (orange, resp. s) or can be "conflicting" (red, resp. c).
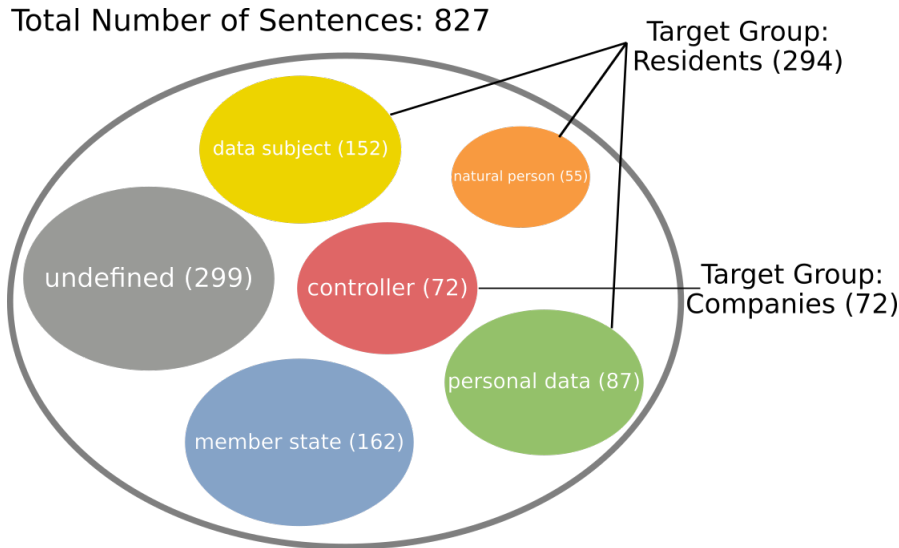
Figure 2: Grouping GDPR Constraints Along Target Groups

# 2 Application to GDPR

A currently very important regulation is the GDPR. This legislation consist of 88 pages and is therefore quite extensive; however, citizens should know their rights regarding data privacy. Besides citizens, companies might also be interested in their duties regarding data privacy of customers. But not every paragraph is equally important for these distinct target groups. For example, the GDPR contains instructions how member states have to enact and adapt this law. These parts of the GDPR might be less relevant for companies or citizens. When applying the presented method, a filtering and grouping based on topics is possible which enables the direct detection of relevant passages for each target group. In this case, we have performed a grouping based on the words "member state" (162), "natural person" (55), "data subject" (152), "personal data" (87) and "controller" (72), inspired by the roles mentioned in [4]. The number in brackets corresponds to the number of sentences we received. An overview of these results is illustrated by Fig. 2.

299 sentences did not contain any of the given words and were therefore categorized as "undefined". This sums up to 827 sentences. The ones that are relevant for citizens, i.e., those containing the words "natural person", "data subject" and "personal data" are altogether 294 sentences. Including the non-categorized sentences reduces the amount of sentences that need to be read from 827 to just 593, which corresponds to a reduction by approximately 28%. Assume that one is interested in what a company's controller needs to take care of. Without considering "undefined" sentences, only 72 sentences need to be evaluated, resulting in a reduction of 91%. If the "undefined" ones are
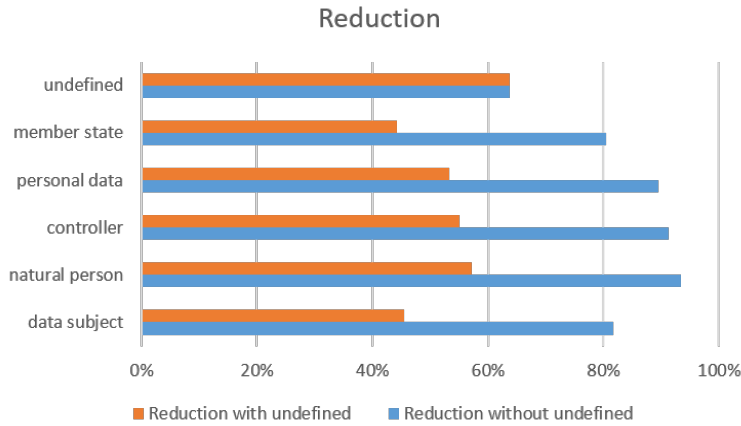
Figure 3: Reductions by Target Group

considered 371 sentences have to be reviewed, still resulting in a reduction of 55%. The highest reduction of 93% can be achieved for target group "natural person" without "undefined " and 57% with "undefined". Note that taking "undefined" into account for one target group is a maximum assumption in the sense that all "undefined" constraints actually belong to exactly this target group. Figure 3 summarizes all reduction times.

To explain this in more detail consider the following sentences from the GDPR:

- *Where personal data are processed for the purposes of direct marketing, the data subject should have the right to object to such processing, including profiling to the extent that it is related to such direct marketing, whether with regard to initial or further processing, at any time and free of charge.*

- *Where personal data are processed for direct marketing purposes, the data subject shall have the right to object at any time to processing of personal data concerning him or her for such marketing, which includes profiling to the extent that it is related to such direct marketing.*

These are displayed as redundant sentences in the output graph and can therefore be handled at once. Reading the document in a chronological order, the first sentence would be on page 13, the second on page 45. It might be difficult to recognize that these are redundant sentences.

Grouping the GDPR as explained before results in the graph displayed schematically in Fig. 4.

Determining the relations between the constraints does not directly reduce the number of sentences that need to be read, but facilitates the implementation and maintenance of regulatory packages such as the GDPR.
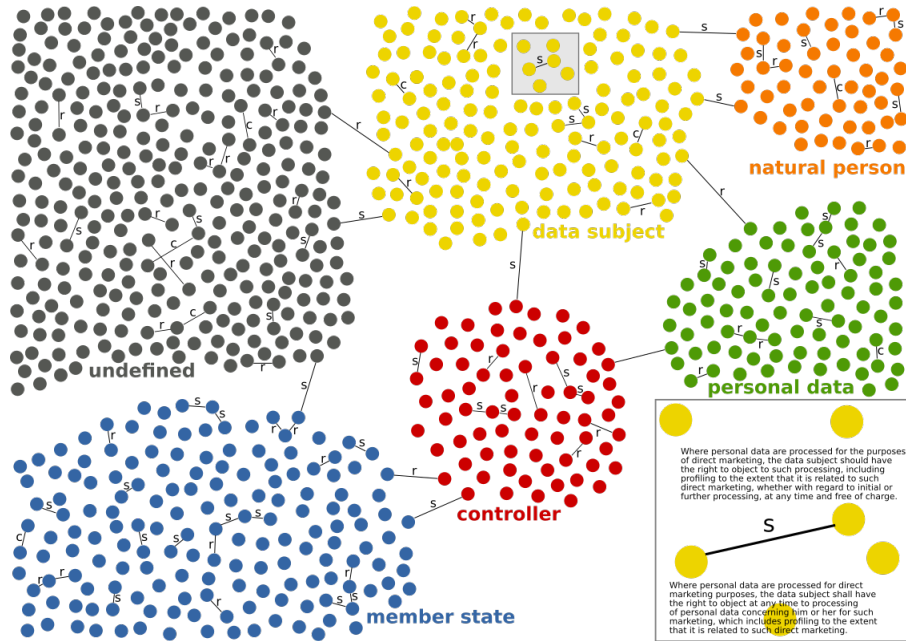
4

Figure 4: GDPR Analysis: Constraints and their Relations

# 3 Related Work

Extracting knowledge from text is a broad and highly regarded task in science and practice. One distinction is the type of text or documents that is analyzed. It ranges from social media texts, e.g., [1] over regulatory documents, e.g., [3, 6, 11, 12] and business process descriptions, e.g., [7, 8, 10] to historic text analysis as in digital humanities, e.g., [5].

Only few approaches target at digitalizing the GPDR such as [4]: here the GDPR is formalized in terms of a declarative notion, the so called DCR graphs.

# 4 Future Challenges

For future work we target to provide our tool as a web service. The input are the regulatory documents. During the application of the ConRelMiner the user can set different parameters for pre-processing and choose between different methods for the processing. This is particularly helpful if users already have some (domain) knowledge about the regulatory documents at hand (e.g., which additional information can be used). However, the ConRelMiner can be also applied without any prior knowledge and without any interaction: just input some documents and receive the filtered and grouped set of constraints, together with their relations.

# Acknowledgment

# References

[1] C. C. Aggarwal and H. Wang. Text mining in social networks. In *Social Network Data Analytics*, pages 353–378. 2011.

[2] C. C. Aggarwal and C. Zhai. *Mining text data*. Springer Science & Business Media, 2012.

[3] I. S. Bajwa, M. G. Lee, and B. Bordbar. SBVR business rules generation from natural language specification. In *AAAI spring symposium: AI for business agility*, pages 2–8, 2011.

[4] S. Debois, T. T. Hildebrandt, P. H. Laursen, and K. R. Ulrik. Declarative process mining for DCR graphs. In *Proceedings of the Symposium on Applied Computing, SAC 2017, Marrakech, Morocco, April 3-7, 2017*, pages 759–764, 2017.

[5] A. Don, E. Zheleva, M. Gregory, S. Tarkan, L. Auvil, T. E. Clement, B. Shneiderman, and C. Plaisant. Discovering interesting usage patterns in text collections: integrating text mining with visualization. In *Proceedings of the Sixteenth ACM Conference on Information and Knowledge Management, CIKM 2007, Lisbon, Portugal, November 6-10, 2007*, pages 213–222, 2007.

[6] M. Dragoni, S. Villata, W. Rizzi, and G. Governatori. Combining nlp approaches for rule extraction from legal documents. In *1st Workshop on MIning and REasoning with Legal texts (MIREL 2016)*, 2016.

[7] F. Friedrich, J. Mendling, and F. Puhlmann. Process model generation from natural language text. In *International Conference on Advanced Information Systems Engineering*, pages 482–496. Springer, 2011.

[8] A. Ghose, G. Koliadis, and A. Chueng. Process discovery from model and text artefacts. In *Services, 2007 IEEE Congress on*, pages 167–174. IEEE, 2007.

[9] F. Nazir, W. H. Butt, M. W. Anwar, and M. A. K. Khattak. The applications of natural language processing (NLP) for software requirement engineering-a systematic literature review. In *International Conference on Information Science and Applications*, pages 485–493. Springer, 2017.

[10] M. Riefer, S. F. Ternis, and T. Thaler. Mining process models from natural language text: A state-of-the-art analysis. *Multikonferenz Wirtschaftsinformatik (MKWI-16), March*, pages 9–11, 2016.

[11] K. Winter and S. Rinderle-Ma. Detecting constraints and their relations from regulatory documents using nlp techniques. In *Int'l Conference on Cooperative Systems*, 2018. (accepted for publication).

[12] K. Winter, S. Rinderle-Ma, W. Grossmann, I. Feinerer, and Z. Ma. Characterizing regulatory documents and guidelines based on text mining. In *On the Move to Meaningful Internet Systems. OTM 2017 Conferences - Confederated International Conferences: CoopIS, C&TC, and ODBASE 2017, Rhodes, Greece, October 23-27, 2017, Proceedings, Part I*, pages 3–20, 2017.