# SEMANTIC TEXT MATCHING OF CONTRACT CLAUSES AND LEGAL COMMENTS IN TENANCY LAW

Jörg Landthaler / Elena Scepankova / Ingo Glaser / Hans Lecker / Florian Matthes

Research Associates, Technical University of Munich, Department of Informatics, Software Engineering for Business Information Systems, Boltzmannstraße 3, 85748 Garching bei München, DE, joerg.landthaler@tum.de, elena.scepankova@tum.de, ingo.glaser@tum.de; http://wwwmatthes.in.tum.de

Rechtsanwalt, Produktmanager der Zielgruppe Rechtsanwälte, Haufe-Lexware GmbH & Co. KG, Fraunhoferstr. 5, 82152 Planegg/München, DE, hans.lecker@haufe-lexware.com; http://www.haufe-lexware.com

Professor, Technical University of Munich, Department of Informatics, Software Engineering for Business Information Systems, Boltzmannstraße 3, 85748 Garching bei München, DE, matthes@in.tum.de; http://wwwmatthes.in.tum.de

***Abstract:*** *We present an innovative approach to support lawyers in the processes of contract drafting, editing and analysing. Legal comments on German tenancy law contain condensed knowledge that supports contract writers. Our approach supports lawyers by providing relevant information from legal comments and practitioner books. Our approach allows to select an arbitrary length text passage of interest and displays relevant corresponding information inline. It is based on text relatedness measures, in particular tf-idf and word embeddings, and is an instance of the semantic text matching problem.*
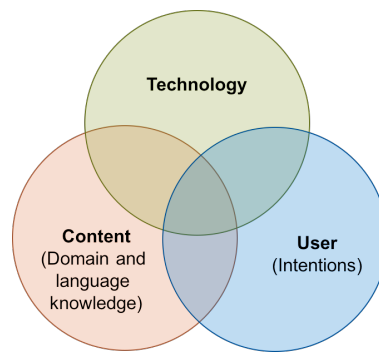
## 1. Introduction

Legal content and media providers face the challenge to support their customers, mostly lawyers, in their work by providing relevant content. However, the relevance of different information depends heavily on the intent of a client and it is a big challenge to identify user intents through an internet application. This paper investigates an approach to support lawyers in the process of contract drafting, editing and analyzing where users actively select a focus text passage of interest and our approach is an attempt to deliver relevant and tailored information on this focus text passage with an unsupervised approach. For example, a clause in a residential lease agreement about cosmetic repairs is subject to a large amount of jurisdiction. Legal comments summarize the existing legislation and judgements on topics relevant for residential lease agreements. This endeavor is subject to different challenges: technological, domain-knowledge and user intention estimation (cf. Figure 1).

Our contribution encompasses an envisioned, innovative approach that delivers relevant information from legal comments for user-selected text passages during the process of contract drafting, editing and analysis to accelerate the involved processes. We show that the problem at hand is an instance of the semantic text matching (STM) problem and propose a purely unsupervised solution. In this paper, we mainly evaluate the technological feasibility of our approach on a dataset of 6 publicly available template tenancy contracts (for private premises) and 3 different legal comments/practitioner books provided by the Haufe Group.

The remainder of this paper is organized as follows: In Section 2, we analyze the process of contract drafting and analysis. Section 3 introduces the general problem of semantic text matching (STM) and proposes a generic solution approach to the STM problem. We show that recommending legal comment chapters for individual contract clauses is an instance of the STM problem. In Section 4, we describe the dataset that we subsequently use in Section 5 to evaluate our approach. Section 6 covers related work, Section 7 concludes the

paper by summarizing related work and conclusions we draw from our work as well as presenting future work.



**Figure 1: Triangle of challenges: to apply our approach in a real-word setting, expertise in different areas is required: The technology used and especially its behaviour; The content at hand, in particular domain knowledge and domain-specific linguistic knowledge; The estimation or elicitation of user intentions.**

## 2.   Process Analysis and Expected Implications

From a company perspective, the drafting, modification and analysis of contracts are part of the contract management. Contract management as a whole is an extensive task for companies which can involve all business units. A major task for in-house attorneys and lawyers is the translation of the contracting parties' intents into legally precise worded contracts in a way that the legal consequences become effective and the contract is legally binding. A part of this process consists of checking whether the contract complies with all requirements (internal and external). Additionally, the task can involve a risk assessment.

Drafting a new contract in the frame of contract negotiations turns very rarely out to be a linear process. In fact, it is more common that the terms of the contract remain subject to modifications until the contract is accepted by all contracting parties. In certain cases, new contracts are based on pre-defined templates which significantly reduce room for negotiations. So-called contract generators can be successfully used in such cases. Long-term contracts usually imply that the contract must be modified over time. Such a modification might be due to external factors, e.g. laws and jurisdiction, as well as changing requirements of the contracting parties. The contract analysis process is based on an existing contract which has to be analyzed from certain perspectives, where both, economic problems and legal consequences have to be taken into account.

In-house attorney's and lawyer's work typically is a knowledge-based process which can be supported by the relevant expertise. Nevertheless, the proposed approach does not stipulate a pre-defined order of steps. In fact, the relevant context can be identified at any point of time and supported by pertinent expertise. Thus, a tool implementing our envisioned approach can be a useful support for the user without restricting the user in his/her method. The approach can be used for drafting, modifying and analyzing contracts since it allows to assess the relevant context based on the text to be considered. Compared to a conventional search, the approach is much more efficient. It is no longer necessary to enter multiple manually phrased search queries because the relevant expertise can be identified easily. Furthermore, the approach can result in a higher quality of expertise support since the user of the content pool does not need any specific skills when entering search queries.
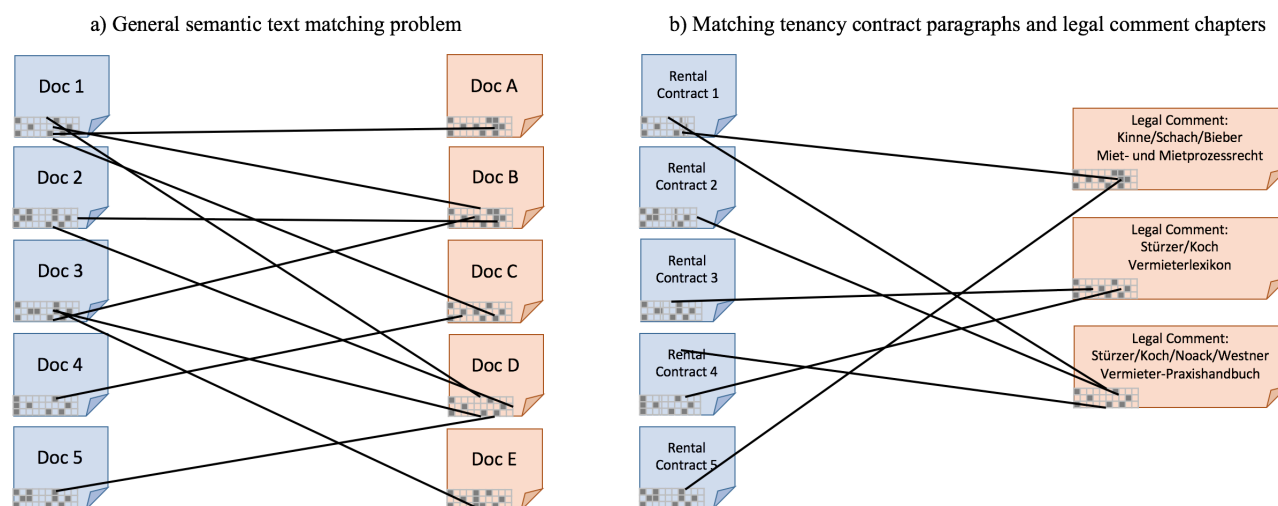
## 3.   Semantic Text Matching and an Application in German Tenancy Law

We introduce the problem of semantic text matching (STM). We define the problem as the identification of logical or semantic links between text fragments of one or several (different) document types. The problem is relevant for many domains, whenever certain documents need to comply with rules laid down in textual form in other documents, cf. Figure 2a). An example that we investigate in this paper is to match contract clauses with legal comment chapters, see Figure 2b). The STM problem reveals different characteristics when examined from different perspectives:

- **(Implicit) References View:** The STM problem typically attempts to identify links between text passages where no explicit link is given in textual form, but the text passages are still logically or semantically linked, cf. [LANDTHALER ET AL. 2016a] and [WALTL ET AL. 2016].

- **Network View:** The STM problem for two different document types can be described using a bi-partite graph as shown in Figure 1 and reveals the assignment/matching character of STM.

- **Search View:** The STM problem is closely related to the general information retrieval (search) problem. However, in contrast to the general information retrieval problem, STM is constrained to specific combinations of document types. This simplifies the general search problem, because by choosing a suitable combination of document types by humans, the existence of actual links is more likely (in contrast to obvious non-relevant results present in general information retrieval). Moreover, text passages are much smaller than full documents, which is typical for information retrieval. The STM also allows to investigate intended and non-intended results in a deeper fashion with less side-effects.

The STM problem can be decomposed into two sub-problems:

- **Segmentation:** Documents of either type (often) need to be segmented into smaller text fragment.

- **Matching:** Text fragments of one document type need to be matched with text fragments of another document type.



**Figure 2: Illustration of Semantic Text Matching (STM) and our application in German tenancy law.**

Sometimes the segmentation problem is already solved for some document types, for example, the legal comment chapters in our application. Otherwise a NLP algorithm or some heuristic needs to be chosen, for example to split a document by markers in semi-structured documents or to apply a sentence segmentation algorithm.

We propose to use text relatedness algorithms to solve the matching problem. While all algorithms for text relatedness / text similarity can be used, we choose the classical TF-IDF measures as a baseline and word embeddings as a competing method. Word embeddings recently gained a lot of attention in the NLP community and in contrast to sparse vector space models that encode the occurrence frequency of individual words / tokens, word embeddings encode the co-occurrence frequency of words / tokens. Word embeddings are vectors that represent individual words. [MIKOLOV ET AL. 2013a] presented a popular method to efficiently calculate word embeddings (trained with a shallow artificial neural network) and showed that word embeddings encode to some degree semantical aspects of words [MIKOLOV ET AL. 2013b], for example synonymy.

## 4. Dataset

We created a dataset in German language consisting of 6 publicly available contract templates (for private premises) and 3 legal comments/practitioner books. The contracts have been obtained in PDF format. We

manually extracted the information to ensure a high quality and manually segmented the clauses into semantically enclosed paragraphs, Table 3 displays examples. Certain legal aspects are present in many or all contracts and usually are slightly differently phrased (with synonyms or re-ordered). Other legal aspects are present in one contract only, for example regarding the topic "energy pass".

| Contract | #Paragraphs | #Words | ∅ Words/Paragraph | #Tags |
|---|---|---|---|---|
| 1 | 146 | 5354 | 31.71 | 195 |
| 2 | 80 | 2267 | 30.85 | 120 |
| 3 | 75 | 2234 | 31.18 | 115 |
| 4 | 90 | 2609 | 30.73 | 134 |
| 5 | 53 | 1304 | 30.40 | 73 |
| 6 | 104 | 3618 | 31.71 | 155 |
| All Contracts | 548 | 17386 | 31.09 | 792 |

**Table 1: Tenancy contract statistics.**

The legal comments have been obtained in an epup/xhtml format and we use the pre-segmented (sub-)chapters. It is worth to note that the intended audience for the legal comment differs. Documents 1 and 2 are intended for legal professionals and contain many references to judgements and cultivate a language that is specifically intended for legal professionals. Document 3 is mostly intended for professional landlords and practitioners. It is shorter in terms of overall size and chapter size, and the language is easier to read for laymen.

| Doc. | Name | #Chapters | #Words | ∅Words/ Chap. | #Tags | # Tags Broad | #Tags Narrow |
|---|---|---|---|---|---|---|---|
| 1 | Miet- und Prozessrecht | 734 | 510129 | 283.08 | - | 96 | 49 |
| 2 | Vermieterlexikon | 890 | 462265 | 256.52 | - | 91 | 52 |
| 3 | Vermieter-Praxishandbuch | 178 | 99569 | 55.25 | 539 | 36 | 26 |
| Tot. | - | 1802 | 1071963 | 198,28 | - | 223 | 127 |

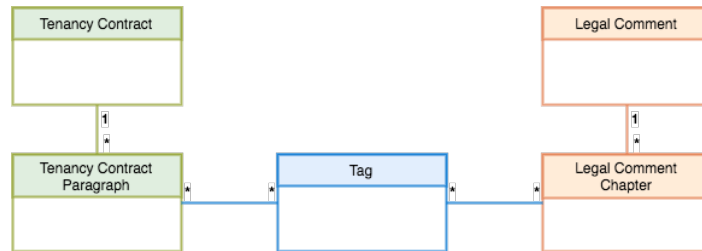**Table 2: Legal Comments statistics.**

The regulations for German template tenancy contracts are complex. Cosmetic repairs are by default a duty for the landlord. The landlord can pass the duty to carry out cosmetic repairs to the tenant. For the effectiveness of passing the duty to the tenant, many different details need to be considered. Slight variations of phrasings easily lead to an ineffectiveness of passing the cosmetic repairs to the tenant. Moreover, certain combinations of individually effective clauses may lead to ineffectiveness due to accumulation effects. Many of these details are described in the legal comments. However, the information is typically spread over several chapters.

| Contract | Paragraph Text | Tags |
|---|---|---|
| 4 | Der (Die) Vermieter __ wohnhaft in __ und der (die) Mieter __ schließen folgenden Mietvertrag:<br>(The landlord __ residing in __ and the tenant __ conclude the following tenancy contract:) | -Contracting Parties |
| 2 | Die regelmäßigen Schönheitsreparaturen während der Mietzeit übernimmt auf eigene Kosten der Mieter.<br>(The regular cosmetic repairs during the rental period have to be done by the tenant on his own cost.) | -Cosmetic Repairs: Tenant Duty |

| | 3 | Während der Mietzeit verpflichtet sich der Mieter, auf seine Kosten erforderliche Schönheitsreparaturen fach- und sachgerecht durchzuführen. (During the tenancy period the tenant feels obliged to do the cosmetic repairs on his own cost in a proper and appropriate way.) | -Cosmetic Repairs: Tenant Duty<br>-Cosmetic Repairs: Kind and Quality |
|---|---|---|---|
| | 1 | Die Schönheitsreparaturen müssen fachgerecht ausgeführt werden. (The cosmetic repairs have to be carried out in a proper way.) | -Cosmetic Repairs: Kind and Quality |

**Table 3: Exemplary text paragraphs from tenancy contracts and the tags we attached to these paragraphs.**

We created different ground truths of links among the contract paragraphs and the legal comment chapters. In order to scan each contract and legal comment only once, we used tags to simplify the process. We started with the contracts and identified regulatory content that is present in at least two of the contracts. The contracts have been fully tagged yielding 214 different labels. One paragraph can hold several tags, because different regulatory contents are spread over two sentences in one contract but covered in once sentence by another contract. The data model of our dataset is illustrated in Figure 3. Regarding the legal comments, we tagged document 3 with respect to all 214 tags. All comments have been tagged regarding nine cosmetic repair tags.



**Figure 3: The data model of tenancy contracts, tags and legal comment paragraphs.**

For tagging the legal comments with respect to the nine tags, we started with a case-insensitive keyword based search for the term "schönheitsreparatur" ("cosmetic repair") that resulted in 159 chapters. We manually tagged the chapters in two ways: a broad and a narrow fashion. The narrow tagging includes chapters where a phrasing example for the phrasing of the particular tag is present, the full chapter covers the topic of the tag or the chapter contains single lines that cover highly relevant aspects for the particular tag (and the respective contract clauses). In a broad tagging, we additionally included chapters that contain indirectly or less important aspects for a particular tag. For example, in the broad tagging we included the comment chapter that informs about the default scope of cosmetic repairs for the tag "passing of the duty". If a "passing of the duty" clause is included in a contract, it might be relevant to know about the implications of not specifying the scope of this duty. In total, this leads to 2046 relations between contracts and legal comment document 3, and 981 and 554 true relations between contracts and all legal comments for the cosmetic repairs in the broad and narrow fashion.

## 5. Evaluation

Our proposed envisioned approach allows the user to select an arbitrary amount of text in a tenancy contract. However, the pre-segmented dataset described in Section 4 enables us to quantitatively evaluate our approach regarding the suitability of different text relatedness measures, pre-processing technologies, training corpuses and segmentation technologies. It also enables us to deeper investigate the linguistic challenges of the content.

We preprocessed all texts by removing xml tags, line breaks, lower-casing and removing all non-alpha characters to obtain clean tokens (SP). In order to retain the paragraph sign, we replaced it with a special token. For all experiments, we iterate through all contract paragraphs and calculate a ranking of all legal comment chapters for each clause using either the tf-idf or the word embeddings text representations and employ the

cosine similarity measure. We measure the quality of the ranking with the RP-Score (RPS), [LANDTHALER ET AL. 2016b], that calculates the average position of legal comment chapters with respect to our ground truth. Additionally, we calculate the average first ranking position (FRP) of the first legal comment chapter that has a true relation to one clause according to our ground truth. Smaller RPS and FRP values indicate better results. A user will likely be presented only a selection of the results, for example the first 10 results, which should be evaluated with for example precision/recall measures. However, the RPS allows us to investigate the reasons for poorly ranked legal comment chapters.

|  |  | Legal Comm. 3 | | CR Broad | | CR Narrow | |
|---|---|---|---|---|---|---|---|
|  |  | FRP | RPS | FRP | RPS | FRP | RPS |
| TF-IDF | Standard preprocessing (SP) | **5.44** | 28.06 | 2.40 | 390.68 | 4.66 | 284.49 |
|  | SP + Stopwords Removed (SR) | 5.64 | 27.37 | **2.43** | 378.72 | 4.43 | 259.51 |
|  | SP + SR + Stemming | 5.49 | **24.96** | 2.93 | 379.16 | 4.06 | 265.63 |
|  | SP + Only POS-tagged nouns (=POSN) | 7.21 | 28.74 | 2.76 | 369.62 | 5.43 | 254.74 |
|  | SP + SENT | 23.41 | 75.50 | 96.80 | 799.99 | 230.23 | 811.94 |
| Word Embeddings | Trained on Contracts + Comment Doc. 3 | 16.97 | 42.49 | - | - | - | - |
|  | Train. on Contracts + All Comm. (CC) | 17.39 | 41.88 | 5.61 | 473.78 | 5.76 | 328.36 |
|  | Train. on CC + GCC | 17.70 | 42.23 | 6.09 | 475.69 | 6.42 | 336.97 |
|  | Train. on CC + GCC + Wikipedia | 19.85 | 47.19 | 28.28 | 599.28 | 33.33 | 474.61 |
|  | Train. on CC + SR ⊗ | 18.36 | 42.28 | 6.04 | 482.18 | 6.85 | 326.09 |
|  | Train. on CC + SR + Stemming(=BEST) | 16.38 | 37.37 | 6.28 | 442.64 | 6.76 | 302.84 |
|  | Train. on CC + POSN ⊗ | 16.47 | 38.67 | 5.47 | 415.07 | 6.23 | 264.45 |
|  | Train. on CC + SENT | 25.85 | 58.17 | 3.25 | **255.33** | **3.33** | **233.29** |

**Table 4: Quantitative evaluation of different technologies, pre-processing methods and training corpuses with RPS and FRP evaluation measures on three different evaluation sets: ⊗ word embeddings trained without stemming or stop words removal. Smaller RPS and FRP values indicate better results. Best results marked in bold.**

We calculate the tf-idf representation of the corpus with genism[1]. Applying pre-processing methods like stopwords removal (SR), stemming and selecting only nouns (Penn treebank tags NN, NNS, NNP, NNPS, calculated with spacy[2]). It is known that nouns carry the semantic information of the subject that a text covers. All pre-processing methods and their combination slightly improve the results for the large corpus of all legal comments. For the smaller corpus consisting of legal comment 3 only, the pre-processing improves the results. We also experimented with the word embeddings text representation. We obtain a vector that represents a clause or a legal comment chapter by summation over all word vectors in a paragraph or chapter. Word Embeddings need to be pre-trained. We use the original word2vec[3] implementation. Commonly, the rule of thumb is that larger training corpuses lead to "better" word embeddings. This seems not to apply for our experiments. We trained word embeddings on different corpuses with 100 iterations and a vector size of 300 and standard parameters for the other word2vec parameters. All of them include the contracts. We trained word embeddings on legal comment 3 and all legal comments. Next, we added the German Civil Code (a relevant law for the topic). We then added a 1 GB selection of the German Wikipedia. From the results, we conclude that a big

---

[1] https://radimrehurek.com/gensim/, version 3.1.0, last accessed January 2018
[2] https://spacy.io/, version 2.0.5, with German core language model, last accessed January 2018
[3] https://github.com/kzhai/word2vec, version 0.1c (for Mac OS X), last accessed January 2018

training corpus does not always result in better word embeddings. Adding legal comments does improve the results, while adding a relevant law does not improve the results.

We applied different pre-processing technologies to the word embeddings. Summation over many vectors may incorporate lots of noise and including only relevant words improves the results for this task, supported by the good results obtained from including word vectors only for nouns or stemming the words before training word embeddings. Additionally, we further segmented the legal comment chapters into sentences using spacy. Spacy's sentence segmentation stumbles upon references and we post-processed it by adding few character "sentences" to the larger sentences. We calculate vectors that represent a sentence from the words in the sentences, rank them using cosine similarity for each contract paragraph, map the sentences back to the legal comment chapter they stem from and select only unique elements (SENT). This leads to another ranking that compares more equal sized text fragments. The additional sentence segmentation leads to very good average results for cosmetic repairs in the larger corpus (all legal comments), but not for the smaller corpus. A deeper investigation of the reasons for this behavior will be necessary in the future. Note that for the SENT approach, the word embeddings based solution is computationally much more efficient since the vector size can be chosen manually. Table 5 shows that the results also vary a lot for the different legal comments.

| Legal Comment | Broad | | | Narrow | | |
|---|---|---|---|---|---|---|
| | Min Rank | Avg. Rank | Max Rank | Min Rank | Avg. Rank | Max Rank |
| 1 | 1 | 13.91 | 1153 | 1 | 26.8 | 1094 |
| 2 | 1 | 29.32 | 1556 | 1 | 47.62 | 1556 |
| 3 | 1 | 24.33 | 1251 | 1 | 36.19 | 1251 |

Table 5: Evaluation of the results for the cosmetic repairs for the different legal comments (SENT).

From a manual inspection of the results for the contracting parties' clauses and the more detailed results for cosmetic repairs (cf. Table 6), we assume that the biggest challenges for the unsupervised approach are to distinguish between very similarly phrased legal comment chapters where some are relevant and others are not, but also to detect only few relevant lines in large legal comment chapters.

| Tag Name | Instances | Min Rank | Avg. Rank | Max Rank | Top 5 | Top 10 | Top 15 | Top 20 | Top 25 | Top 50 | Rest |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Tenant Duty | 185 | 1 | 268.71 | 1094 | 9 | 14 | 6 | 6 | 4 | 20 | 126 |
| Periods Start | 10 | 3 | 45.30 | 103 | 2 | 0 | 0 | 2 | 1 | 0 | 5 |

Table 6: Detailed results for selected cosmetic repair tags (for the narrow selection: 195 true links).

## 6. Related Work

In [WALTL ET AL. 2016], we identified different implicit references between paragraphs in legal documents. While certain legal documents explicitly hint to other textual passages, many logical or semantical links are not explicitly documented, for example special terms are used without linking to the place where the terms are legally defined. We attempted to extract implicit links between paragraphs in the German Civil Code by counting shared nouns between paragraphs [LANDTHALER ET AL. 2016a]. Our newly coined term "Semantic Text Matching" differs from semantic matching that applies to matching nodes in ontologies [GIUNCHIGLIA ET AL. 2007]. A related technology applied in legal question answering is "textual entailment" that identifies relations between two text entities where one text entity proves a hypothesis present in the other text entity. The STM problem is present in many domains and applications, for example in argumentation mining. The STM problem here is to match text fragments that contain premises or claims, and conclusions or statements, see for example [NADERI & HIRST 2016], [RINOTT ET AL. 2015]. They use word embeddings to calculate the relatedness of text fragments by calculating paragraph vectors like our approach or by comparing word

vectors pair-wise. However, to the best of our knowledge, the STM has not been identified as an individual (re-occurring) problem in natural language processing or otherwise.

## 7. Conclusion

We presented a vision for an inline search that supports contract drafting and analysis by displaying information from legal comments in the domain of tenancy law for arbitrary text passages in contracts. We identified three major challenges for this endeavour: technology understanding, content knowledge and user intentions. We showed that the problem at hand is an instance of the novel semantic text matching problem and presented an unsupervised solution for the problem based on text relatedness measures. We created a dataset including a ground truth to evaluate our approach mainly from a technical point of view (with rather general assumptions on the user intentions). We can report promising results that our vision can be realized to some degree. Key results are that word embeddings do no outperform classical tf-idf measure in general. However, in this scenario word embeddings can lead to better results when the legal comment chapters are further segmented into sentences. For this specific domain, the "larger training corpuses lead to better word embeddings" rule of thumb seems not to apply. Our evaluation mainly tackles the technological challenges. Currently, we develop a prototypical implementation that will enable us to evaluate the approach with real world users. Further research on user intentions during the processes of contract drafting and analysis is required. From a content and technology point of view, the influencing factors for relevant and irrelevant returned legal comment chapters need to be investigated in more depth, as well as a deeper investigation which technology performs better in certain cases. Finally, it will be necessary to find solutions to detect relevant chapters that are hard to find for simple text relatedness measures and to distinguish very similar legal comment chapters.

## 8. References

DAGAN, IDO/ DOLAN, BILL/ MAGNINI, BERNARDO/ ROTH, DAN, Recognizing textual entailment: Rational, evaluation and approaches, Journal of Natural Language Engineering, Cambridge University Press, 2009

GIUNCHIGLIA, FAUSTO/ YATSKEVICH, MIKALAI/ SHVAIKO, PAVEL, Semantic Matching: Algorithms and Implementation, Journal on Data Semantics IX, Springer-Verlag, Berlin, Heidelberg, 2007

KINNE, HARALD/ SCHACH, KLAUS/ BIEBER, HANS-JÜRGEN, Miet- und Mietprozessrecht, 7. Auflage 2013, Haufe Lexware

STÜRZER, RUDOLF/ KOCH, MICHAEL, Vermieterlexikon, 15. Auflage 2017, Haufe Lexware

STÜRZER, RUDOLF/ KOCH, MICHAEL/ NOACK, BIRGIT/ WESTNER, MARTINA, Vermieter-Praxishandbuch, 9. Auflage 2017, Haufe Lexware

LANDTHALER, JÖRG/ WALTL, BERNHARD/ MATTHES, FLORIAN, Unveiling References in Legal Texts - Implicit versus Explicit Network Structures, IRIS: Internationales Rechtsinformatik Symposium, Salzburg, Austria, 2016a

LANDTHALER, JÖRG/ WALTL, BERNHARD/ HOLL, PATRICK/ MATTHES, FLORIAN, Extending Full Text Search for Legal Document Collections using Word Embeddings, Jurix: International Conference on Legal Knowledge and Information Systems, Sofia Antopolis, France, 2016b

MIKOLOV THOMAS/ CHEN, KAI/ CORRADO, GREG/ DEAN JEFFREY, Efficient Estimation of Word Representations in Vector Space, International Conference on Learning Representations, 2013a

MIKOLOV THOMAS/ SUTSKEVER, ILYA/CHEN, KAI/ CORRADO, GREG/ DEAN JEFFREY, Distributed Representations of Words and Phrases and their Compositionality, NIPS'13: Proceedings of the 26th International Conference on Neural Information Processing Systems, 2013b

NADERI, NONA/ HIRST, GRAEME, Argumentation Mining in Parliamentary Discourse, Principles and Practice of Multi-Agent Systems: International Workshops: IWEC 2014, Gold Coast, QLD, Australia and Bertinoro, Italy, October 26, 2015, Revised Selected Papers", Springer International Publishing, Cham, Germany, 2016

RINOTT, RUTY/ KHAPRA, MITESH M./ ALZATE, Carlos/ DANKIN, LENA/ AHARONI, EHUD/ SLONIM, NOAM, Show me your evidence - an automatic method for context dependent evidence detection, Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, EMNLP 2015, Lisbon

WALTL, BERNHARD/ LANDTHALER, JÖRG/ MATTHES, FLORIAN, Differentiation and Empirical Analysis of Reference Types in Legal Documents, Jurix: International Conference on Legal Knowledge and Information Systems, Sofia Antopolis, France, 2016