

Threshold-Based Retrieval and Textual Entailment Detection on Legal Bar Exam Questions

Sabine Wehnert

sabine.wehnert@ovgu.de

Otto von Guericke University Magdeburg
Germany

Wolfram Fenske

wolfram.fenske@ovgu.de

Otto von Guericke University Magdeburg
Germany

Sayed Anisul Hoque

sayed.hoque@st.ovgu.de

Otto von Guericke University Magdeburg
Germany

Gunter Saake

saake@iti.cs.uni-magdeburg.de

Otto von Guericke University Magdeburg
Germany

ABSTRACT

Getting an overview over the legal domain has become challenging, especially in a broad, international context. Legal question answering systems have the potential to alleviate this task by automatically retrieving relevant legal texts for a specific statement and checking whether the meaning of the statement can be inferred from the found documents. We investigate a combination of the BM25 scoring method of Elasticsearch with word embeddings trained on English translations of the German and Japanese civil law. For this, we define criteria which select a dynamic number of relevant documents according to threshold scores. Exploiting two deep learning classifiers and their respective prediction bias with a threshold-based answer inclusion criterion has shown to be beneficial for the textual entailment task, when compared to the baseline.

CCS CONCEPTS

• **Information systems** → **Question answering**; *Similarity measures*; *Relevance assessment*; • **Computing methodologies** → Neural networks.

KEYWORDS

legal text retrieval, textual entailment, stacked encoder, explainable artificial intelligence, threshold-based relevance scoring

ACM Reference Format:

Sabine Wehnert, Sayed Anisul Hoque, Wolfram Fenske, and Gunter Saake. 2019. Threshold-Based Retrieval and Textual Entailment Detection on Legal Bar Exam Questions. In *Proceedings of COLIEE 2019 workshop: Competition on Legal Information Extraction/Entailment (COLIEE 2019)*. ACM, New York, NY, USA, 9 pages.

1 INTRODUCTION

Nowadays, globalization poses a challenge for many international organizations, since they need to ensure compliance to laws of all

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

COLIEE 2019, June 21, 2019, Montreal, Quebec

© 2019 Copyright held by the owner/author(s).

jurisdictions falling under the scope of their activities. Tracking changes in law is a challenging task, especially in statutory law where a single modification may affect the applicability of several legal articles, due to implicit co-dependencies between these documents. While domain experts are mostly required to ensure a reliable assessment of relationships among laws and their implications, the amount of legal documents is hard to oversee for a single person. Therefore, a decision support system can help in finding relevant laws and applying them to a specific question or statement.¹ Finding out whether a statement is true, given a corpus of legal text, falls under the task of legal question answering. A legal question answering system consists of two major parts: document retrieval and textual entailment recognition. In the retrieval phase, relevant law articles are selected for a query, having the form of a statement which shall be supported or contradicted by the law articles from the document collection. During the textual entailment phase, the query and accordingly retrieved legal documents are processed by a classification algorithm which returns “yes” in case of positive textual entailment or “no” otherwise. This work is a contribution to the Competition on Legal Information Extraction/Entailment (COLIEE) competition which provides a dataset from Japanese bar exam questions (translated to English) for evaluating the system performance on both tasks, retrieval and entailment classification. Our contribution involves the following methods:

- We combine results from BM25 scoring with word embedding-based retrieval.
- We develop a stacked encoder ensemble for entailment detection.
- We use thresholding for both approaches.

The remainder of this work is structured as follows: Section 2 outlines related work for both tasks with respect to their achievements using similar methods to our approach. In Section 3, we describe basic concepts for string representation in machine learning models, scoring methods and stacked encoders. We explain our approach in detail in Section 4 and show evaluation results in Section 5. After discussing those results, we conclude our findings and mention our future work considerations.

¹The work is supported by Legal Horizon AG, Grant No.:1704/00082

2 RELATED WORK

The related work for our approach is divided in two parts: The legal information retrieval task and the entailment detection task. The first part consists of approaches using BM25 scoring or word embeddings, as well as similarity thresholding for a retrieval task. We further present deep learning methods, followed by approaches using thresholds for a textual entailment task.

2.1 Legal Information Retrieval

2.1.1 BM25-Based Solutions. In the COLIEE '16 competition, Onodera and Yoshioka apply BM25 scoring for information retrieval with several extensions using query keyword expansion. Their best result was an F-measure of 54.5% [11]. Arora et al. observe the best score with the BM25 scoring method on a different task of legal document retrieval [1], compared to language models and term frequency - inverse document frequency (TF-IDF) weighting. This finding contradicts the previous observations from COLIEE competitions and the FIRE 2017 IRLed Track, where ranking SVMs [12] or language models [17, 32] performed better than mere BM25 scoring. Despite those observations, BM25 has shown to provide at least competitive results in many cases, so that we consider it as part of our approach.

2.1.2 Word Embeddings. Word Embeddings have proven to be useful in many natural language processing contexts. We outline several works which have used this document feature representation for legal information retrieval. During the COLIEE '18 competition, the SPABS team was able to overcome vocabulary mismatch in some cases using an RNN-based solution with Word2Vec embeddings trained on English legal documents [34]. Team UB used word embeddings with PL2 term weighting [34]. Yoshioka et al. suggest to use semantic matching techniques for hard questions involving vocabulary mismatch combined with more reliable lexical methods for easy questions [34]. This is the main motivation for our retrieval system, which incorporates lexical BM25 scoring and word embeddings as a semantic representation, respectively.

2.1.3 Thresholding. Thresholding based on similarity values can improve retrieval results by filtering out low-scoring matches. Islam and Inkpen use similarity thresholds to increase the precision of text matching [10]. Stein et al. also employ thresholds for plagiarized document retrieval [29]. In the COLIEE '18 competition, team UBIRLED use a similarity threshold for filtering out irrelevant case judgments [13]. Nanda et al. select the top-5 matching documents from a topic clustering approach [20]. Given the document with the highest similarity score to the query, they apply thresholding, such that any further document will be incorporated into the result set if the distance to the topmost document is less than 15%. Our approach uses a similar criterion for document inclusion.

2.2 Legal Textual Entailment

2.2.1 Deep Learning Approaches. Deep learning approaches have been used by several authors for entailment detection, starting with an application of a single-layered long short-term memory network (LSTM) for input encoding by Bowman et al. [3]. The encoded features from both texts are concatenated and passed through three 200-dimensional tanh layers to a softmax classifier for predicting

the entailment relationship. The task is performed on the SNLI² dataset which is based on image captioning. Rocktäschel et al. apply neural attention [2] for entailment recognition on the same SNLI corpus [26]. Two LSTM networks are employed for encoding the query and the document, whereby the output vectors from the document are used by an attention mechanism for each word in the respective query. Their method achieves 83.5% accuracy, which is compared to the results by Bowman et al. an improvement of 3.3 percentage points. Liu et al. use a bidirectional LSTM with an attention mechanism [16] and obtained 85% accuracy on the SNLI dataset. A stacked encoder architecture developed by Nie and Bansal achieved 86.1% accuracy on the SNLI dataset. Considering that result as the state-of-the-art, we adapt the main idea to our task in the legal domain and further explain this architecture in section 3.3. Do et al. use a convolutional neural network (CNN) with word embeddings [6]. They incorporate additional features from a TF-IDF and latent semantic indexing (LSI) representation of the sentences. Finally, they feed these features in conjunction with the output of the CNN model into a multi-layer perceptron (MLP) network to predict the answer.

We are inspired by the work of the Chen et al., which focuses on a factoid question and answering system [5]. Their goal is to predict a sequence in the document to answer the query, as opposed to our task of detecting an entailment relationship. They trained two multi-layer bi-directional LSTMs to encode the articles and the query. For encoding the article, they extract multiple features from the query and document pairs: word embeddings of the document (300-dimensional Glove embeddings), an exact matching flag, token features (part-of-speech tags, named entity tags, normalized term frequencies) and attention scores for the similarity of a document and the aligned query. These features are concatenated to form the input vector for the LSTM that encodes the article. The question is encoded without extracting any features. Their evaluation is based on the top five pages returned by the algorithm, and results in 77.8% of correct answers on the SQuAD [23] dataset.

Nanda et al. apply a hybrid network of LSTM networks coupled with a CNN, with the final prediction based on a softmax classifier [20]. They use pre-trained general-purpose word embeddings from the Google news corpus, consisting of 3 billion words. Their accuracy for the COLIEE '17 competition was 53.8%, which they attribute to the general-purpose embeddings which may not capture important semantic relationships needed for the legal domain.

From these works, we conclude that LSTM architectures are suitable for entailment detection for open-domain tasks. However, the COLIEE dataset poses a challenge for deep learning models due to the specific meaning of terms in the legal domain and the rather small size of the dataset. Therefore, we refrain from training word embeddings on the statute law competition corpus only, but consider using other general-purpose word embeddings and a slightly different architecture compared to the previous work. We also find in the related work that extracting additional features from the documents can improve the classifier performance.

2.2.2 Thresholding. Thresholding for the entailment task is applied in two cases: First, the entailment detection can be done by using a similarity threshold. This works similar to an attention layer

²nlp.stanford.edu/projects/snli/

in a neural network, which applies a focus on a subsection of the input. The second case refers to thresholding in classifier output probabilities. During the COLIEE '15 competition, Kano obtained good results using a threshold for snippet scoring of the entailment task in the runs KanoLab1, KanoLab2 and KanoLab3 [12]. Ha et al. show that similarity scoring can be used for entailment detection [9], however, this may only work for low-level inferences. Instead, Carvalho et al. suggest to use classifiers with a rich feature set [4]. We incorporate an exact matching component for obvious cases of positive entailment in our solution and train a classifier for the majority of cases. Glickman et al. tune their probabilistic entailment detection method with a threshold [8]. Rooney et al. note that for entailment detection, exploiting individual classifier bias with a voting scheme or stacking is common [27]. We employ a voting scheme which is based on using thresholds on the output probabilities of classifiers in our ensemble and consider their respective bias.

3 BACKGROUND

In this section, we introduce the required concepts to understand the components of our approach. First, we describe two text representation approaches. Second, we present scoring methods for the retrieval task. Finally, we explain a method for sequence encoding and decoding which is frequently applied in open-domain question answering.

3.1 Text Representation

In order to apply machine learning techniques on text data, a suitable representation of the input is needed. We consider two alternative approaches: the bag-of-words model and word embeddings. The bag-of-words representation collects all words of a document regardless of their order as assigns each distinct word to a unique index. For two texts to be considered as similar, there needs to be a significant lexical overlap. In contrast, word embeddings are semantic representations of a word, which do not require a lexical overlap. However, these representations need to be created, for instance from a neural network which is trained on a reference corpus to predict the next probable word in a sequence. A model for training word embeddings is Word2Vec [19]. The Word2Vec algorithm learns the word embeddings by using the continuous bag-of-words (CBOW) model or the continuous skip-gram model. The embeddings are learned in the CBOW model by predicting the current word based on a context window. In contrast, the continuous skip-gram model predicts the surrounding context given a reference word. General-purpose pre-trained word embeddings are often used, for instance Glove embeddings [22]. During the embedding training, the hidden layers of this network capture information such as the co-occurrence of two words which are present in the same context. The resulting n -dimensional weights from the hidden layer of the network are then called word embeddings. It is important in this context to choose a good reference corpus which represents the vocabulary and common word use in the domain of the task well to obtain adequate embedding weights.

3.2 Weighting and Scoring Methods

Regardless of the chosen representation method, there can be the need to adjust the weight of each word while aggregating to the document level. A common approach for weighting is a discounting of the term frequency within a document by its overall frequency within the corpus. This term frequency - inverse document frequency (TF-IDF) weighting scheme takes into account that only certain words describe a document well. That is, words which occur frequently within a document, but relatively seldom across the rest of the document collection are considered to be keywords and obtain the highest weight. Likewise, words which occur often in the whole corpus are treated as stopwords and their weight is diminished by the denominator term of the inverse document frequency. Generally, retrieval systems often contain a ranking function, where the similarity or a relevance score between the query and candidate documents is computed in order to determine document relevance. We describe two scoring methods in more detail: BM25 similarity scoring and Word Centroid Distance, whereby the latter treats the analogous problem of distance measurement for similarity scoring. BM25 scoring, also referred to as Okapi BM25, is similar to TF-IDF weighting, however it limits the influence of the term frequency and common words between query and document, following a notion of term eliteness as a poisson distribution [25]. The Word Centroid Distance (WCD) has been specifically developed for word embeddings by Kusner et al. [15]. It depicts the minimum cumulative distance that the words of a document need to travel to reach all words of another document in the embedding space. Thereby, a semantic distance metric is provided, given that word embeddings exist which capture those semantic relationships between the words of the respective domain. As a result, paraphrasing documents which convey the same meaning without sharing a single word can be still detected as a match.

3.3 Sequence Encoding and Decoding

Using sentences or article paragraphs as an input for a machine learning algorithm requires a transformation to a lower-dimensional feature representation. Since the order of words and thereby also dependencies among words can impact the entailment relationship, we briefly introduce an approach to model text sequences: the sequence-to-sequence (Seq2Seq) model by Sutskever et al. [31]. It consists of a recurrent neural network (RNN) as an encoder, which maps the input to a fixed-length context vector. This context vector is then accessed by the decoder RNN to generate the target sequence. This approach is popular in the general question answering domain [24]. The Seq2Seq model has been extended by Bahdanau et al. by an attention layer, so that a smaller segment of the fixed-length vector can be used by the decoder to predict the target [2]. Neural network architectures built for question answering encode the query and the document separately. A common practice is to concatenate both context vectors to a single context representation and separate them by markers of question start and sequence end for the final prediction [30]. The question answering problem we consider in this work is reduced to a binary output of confirming or rejecting the statement from the query. Therefore, the decoder is replaced by an entailment classifier, as in the work by Nie and Bansal [21]. While they use an MLP layer followed by a

softmax layer for the prediction, we employ a multi-layer neural network, with stepwise-linear activation functions, the rectified linear units (ReLU), due to their accuracy and simplicity at the same time. In particular, recent studies have investigated the approximation capability of ReLUs in hidden layers, and found that there is a gain in accuracy, although the nonlinearity of weights is given up [28, 33]. Schmidt-Hieber names desirable properties of the ReLU function: First, the projection property of ReLUs can be exploited to propagate a signal over layers without distortion to synchronize smaller subnetworks of varying depth, and second, the upper bound of all network parameters equals one, thus limiting the amplitude of weight updates during training [28]. To conclude, sequence encoding has been applied on several question answering problems and therefore, we consider it for the task of legal bar exam entailment classification. In the next section, we present the details of our approach.

4 THRESHOLD-BASED RETRIEVAL AND ENTAILMENT DETECTION

In this section, we present our methodology to retrieve relevant laws for a given query. Then, we show our approach for classifying the entailment relationship.

4.1 Law Retrieval combining BM25 Scoring and Word Embeddings

Previous competition results suggest that it can be beneficial to combine the advantage of lexical matching (using the bag-of-words approach) with more recent semantic matching techniques (e.g., word embeddings) [13].

4.1.1 Conceptual Approach. Our retrieval system uses lexical BM25 scoring together with Word2Vec embeddings. An overview of the process is shown in Figure 1. During the preprocessing phase, we separate the articles in the Japanese civil law file by using regular expressions. Then, we create an index over all articles and apply the BM25 scoring method. In addition, we implement a separate pipeline for word embedding creation from the German civil code as an auxiliary corpus. Preliminary experiments have shown that solely using Word2Vec on the English translation of the Japanese civil code does not provide a good embedding quality, so we enriched the corpus for the embedding creation by an English translation of the German civil code³. With our trained word embeddings, we applied the Word Centroid Distance (WCD) as a scoring function. Further experiments led us to apply TF-IDF weighting on the Word Centroid Distance to increase the number of relevant documents. We find that both methods returned similar results, however, the word embeddings can contribute correct documents in case of vocabulary mismatch and sufficient semantic similarity. Therefore, we introduced similarity thresholding and prioritize the word embedding-based matching documents over the ones which scored high only on BM25. Since the number of relevant documents varies by query, we define criteria for similarity thresholds in both methods to select one or multiple documents, if there is a high similarity to the query. Setting a high similarity threshold results

in a high precision, but a low recall. We therefore select all documents with a high similarity automatically as relevant. Additional documents are only included if they fulfill the threshold criteria which we set manually to optimize retrieval performance on the training data. For instance, a top-1 document with high confidence has either a BM25 score that is at least twice as high as the following second document and larger than the length of the query with an added constant of 20, or that document has a word embedding-based similarity of 90%. The length of the query is considered as a criterion for the BM25 results because the lexical overlap should be significant and not fall below 20% similarity score. Further criteria for lower confidence results are based on the same logic, just with relaxed parameters. Depending on the number of already accepted predictions from higher confidence votes, we define criteria for further result inclusion. For example, if there are no predictions from high, medium or low similarity thresholds, the top-1 prediction without any confidence is retrieved from the results of the word embedding-based similarity scoring.

4.1.2 Implementation Details. Our implementation for the retrieval task is based on Python 3 and Elasticsearch 6.4.1. Furthermore we use the library `vec4ir`⁴ provided by Galke et al. for word embedding-based retrieval and re-weighting frequent terms with the inverse document frequency before the centroids are calculated for the WCD [7]. The 300-dimensional Word2Vec embeddings are trained with the continuous bag-of-words model and a context window size of 5. We train the embeddings for 700 epochs on the English version of the German Civil Code and then for further 800 epochs on the English version of the Japanese Civil Code. The input string is tokenized using the `spaCy`⁵ library.

4.2 Textual Entailment using an ensemble of Stacked LSTM Encoders

4.2.1 Conceptual Approach. Our approach for entailment detection consists of stacked encoders and an entailment classifier, inspired by the stacked bidirectional LSTM encoder by Nie and Bansal [21]. Since there may be identical sequences of query and document, such easy cases are returned as positive entailment beforehand. The remaining query-document pairs are processed by the stacked encoder approach. An overview of our classification process is shown in Figure 2. As a first step, the queries and articles are preprocessed by tokenizing the sequences. Both have a separate pipeline, since four additional features are extracted from the article. The first three features are obtained by comparing the article with query tokens and storing matches as a binary feature if a match occurs in the original form, in the lower-cased form or in the lemmatized token form, respectively. We also compute the normalized term frequency for each term in the article. The main advantage of using the normalized term frequency is to lessen the effect of high occurrences of a term in a text sequence [18]. Once extracted, these four features along with the tokenized form of the article and the query become the input to the deep learning model. The feature vectors of query and article are passed to the stacked encoders separately. The first layer in the pipeline of the deep learning model is the

³https://www.gesetze-im-internet.de/englisch_bgb/

⁴<https://github.com/lgalke/vec4ir/>

⁵<https://spacy.io/>

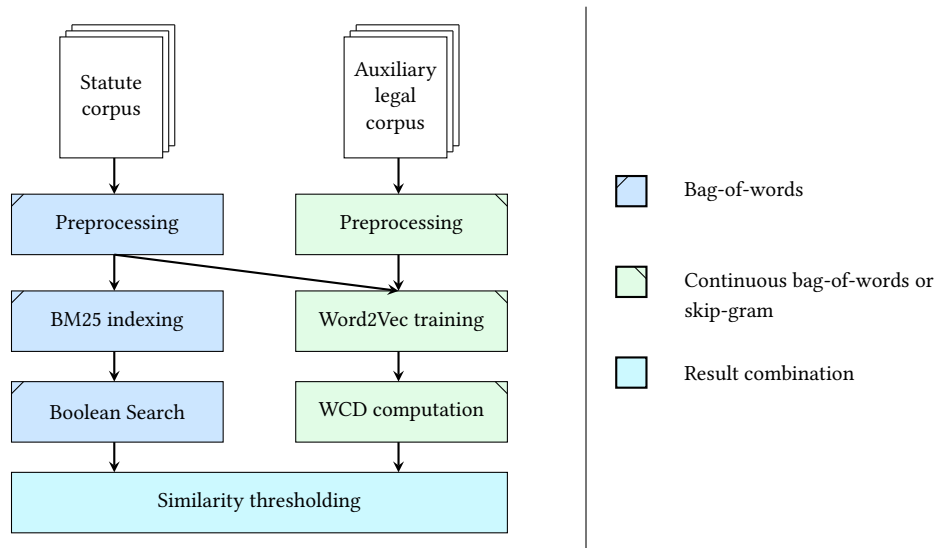


Figure 1: Overview of the information retrieval workflow.

embedding layer. Due to a few spelling mistakes, some words in the documents could not be mapped to the Glove embeddings. Because of the small corpus size, we defined a dictionary with correct term spellings. If the Glove embedding is not found, the correct surrogate word is looked up, at least for cases we detected in the training dataset. We then incorporate an attention layer because of the different size of the query and article length (with the article being potentially longer). Another reason for choosing an attention layer is the interpretability of the learned input characteristics. Words causing higher neuron activation can be traced from the attention scores of the model. As a next layer, this network uses a deep bidirectional long short-term memory (LSTM) network and outputs a fixed-length context representation vector. The context vectors from both the article and the query encoder are concatenated and processed by the entailment classifier consisting of stacked linear layers with ReLU activation functions. The binary output is generated by a sigmoid function providing the predicted entailment relationship. Afterwards, models with the highest accuracy are investigated with respect to their bias to set rules for the ensemble voting mechanism. We use the probability of each class - positive or negative entailment - as an indicator of classifier confidence. For example, if one model predicts positive entailment while the other model votes for negative entailment, the positive result is chosen if the probability of the positive class in the approving model is larger than 53%.

4.2.2 Implementation Details. The entailment recognition task is also implemented using Python 3. For preprocessing we again use spaCy. The deep learning models are built by using PyTorch⁶. We train two models using different random seed values for 100 epochs with a batch size of 20 and an Adamax optimizer with the cross entropy loss. There are 3 bidirectional LSTM layers in both encoders - for query and article features - with a hidden size of 100. We apply

a dropout of 20% for each bidirectional LSTM layer and use gradient clipping. To generate the prediction, there are 4 stacked linear layers with ReLU activation functions, followed by a sigmoid function for the final output.

5 RESULTS

The evaluation results are based on the assessment of the COLIEE competition. First, we present and discuss the results of the retrieval task. Second, we proceed with the results of the textual entailment task and insights from the assessment. While a complete question answering system is evaluated based on the whole pipeline - retrieval and entailment classification on the retrieved documents - we decided to solve both tasks independently from each other by taking the gold standard set of all relevant documents for the entailment recognition.

5.1 Law Retrieval

We summarize our document coverage for high, medium and low confidence in the training dataset in Table 1. The BM25 scoring method achieves a higher coverage than the word embedding-based solution and has therefore a slightly worse retrieval performance on the high, medium and low confidence thresholds. We decide to prioritize word embedding-based results over the BM25 scoring when the similarity - thus also the confidence - is high. For the combined result on the training dataset (H18-H29), we achieve a precision of 29.3%, a recall of 63.59% and an F2-measure of 48.2%. Preliminary experiments have shown that our approach outperforms simple TF-IDF scoring in Elasticsearch on the training data, as well. Table 2 contains our result in the COLIEE competition compared to the competitors with the highest and the lowest score.

Similarly, on the macro-average evaluation setup, we achieve an F2-score of 46.59%, while other participants scored from 40.08% to 54.93%. However, results from the competition runs indicate that our Word2Vec embeddings cannot guarantee relevant document

⁶<https://pytorch.org>

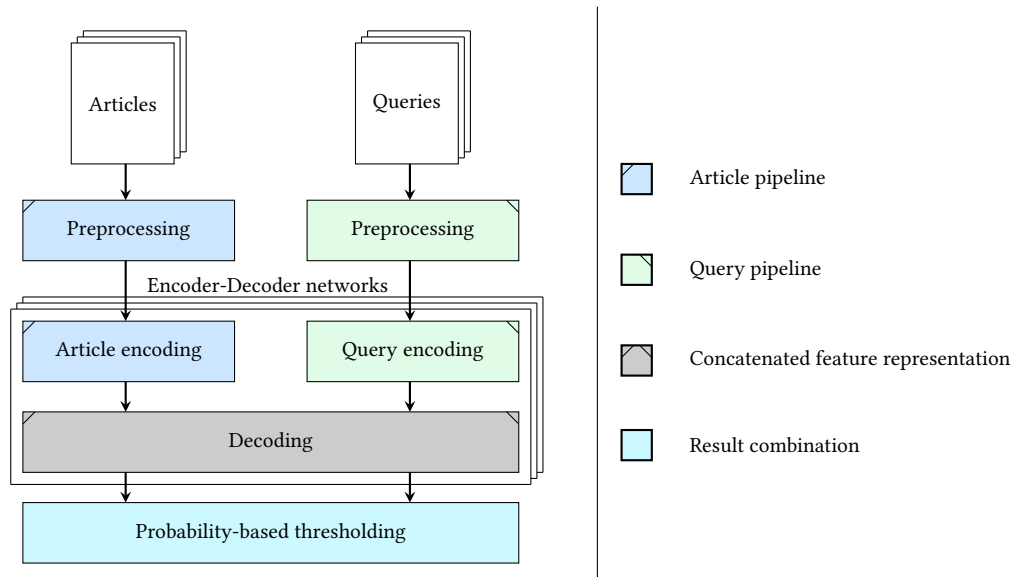


Figure 2: Overview of the textual entailment workflow.

Table 1: Retrieval coverage and performance by using similarity thresholds for several confidence levels. The number of documents retrieved by using BM25-based similarity thresholding is denoted as N_{bm25} , the respective word embedding similarity-based document count as N_{emb} . The metrics coverage C , precision P , recall R and F2-measure $F2$ are depicted in percentages.

Metric	High	Medium	Low
N_{bm25}	89	390	453
N_{emb}	53	154	380
C_{bm25}	8.1	35.8	41.7
C_{emb}	4.9	14.0	31.6
P_{bm25}	97.7	87.5	77.6
P_{emb}	100.0	94.9	81.0
R_{bm25}	94.5	79.5	66.3
R_{emb}	94.3	86.9	78.0
$F2_{bm25}$	94.2	79.5	66.9
$F2_{emb}$	94.8	87.2	76.1

detection, since our recall at 30 has been outperformed by all other methods. The Word2Vec embeddings may be trained on a larger corpus of legal text - not only civil codes - to provide more valuable semantic representations. Preliminary experiments of using the skip-gram model instead of CBOW for training word embeddings decreased all scores on the training dataset, so that we only selected embeddings based on the CBOW model for result submission.

5.2 Textual Entailment

We selected two models with the highest accuracy. They are both biased towards the negative class, although we divided the dataset in such a way that there are more positive examples in the training

Table 2: Retrieval results of our team DBSE for the metrics F2-measure $F2$, precision P , recall R , mean average precision MAP , recall at 5 $R@5$, recall at 10 $R@10$ and recall at 30 $R@30$, calculated as macro-averaged values and depicted in percentages.

Metric	UA-TFIDF	DBSE	iitptfidf
F2	54.93	46.59	40.08
P	59.18	45.44	43.88
R	54.42	49.32	39.63
MAP	61.81	51.19	50.56
$R@5$	61.98	51.24	57.02
$R@10$	69.42	61.98	62.81
$R@30$	76.03	66.94	75.21

dataset (with a validation set from H28). We obtain from both models 63.6% accuracy on the validation set. We consider the models to be complementary because they predict positive entailment on different instances. Therefore, we combine their predictions.

Table 3: Entailment classification results of our team DBSE for the accuracy A metric, calculated as macro-averaged value and depicted in percentages.

Metric	UA_Ex	DBSE	Baseline
A	68.37	57.14	52.04

The assessment from the COLIEE competition shown in Table 3 resulted in an accuracy of our model of 57.14%, compared to the best result of 68.37% and the lowest score of 44.9%. The baseline which always predicts a negative entailment achieved 52.04% accuracy.

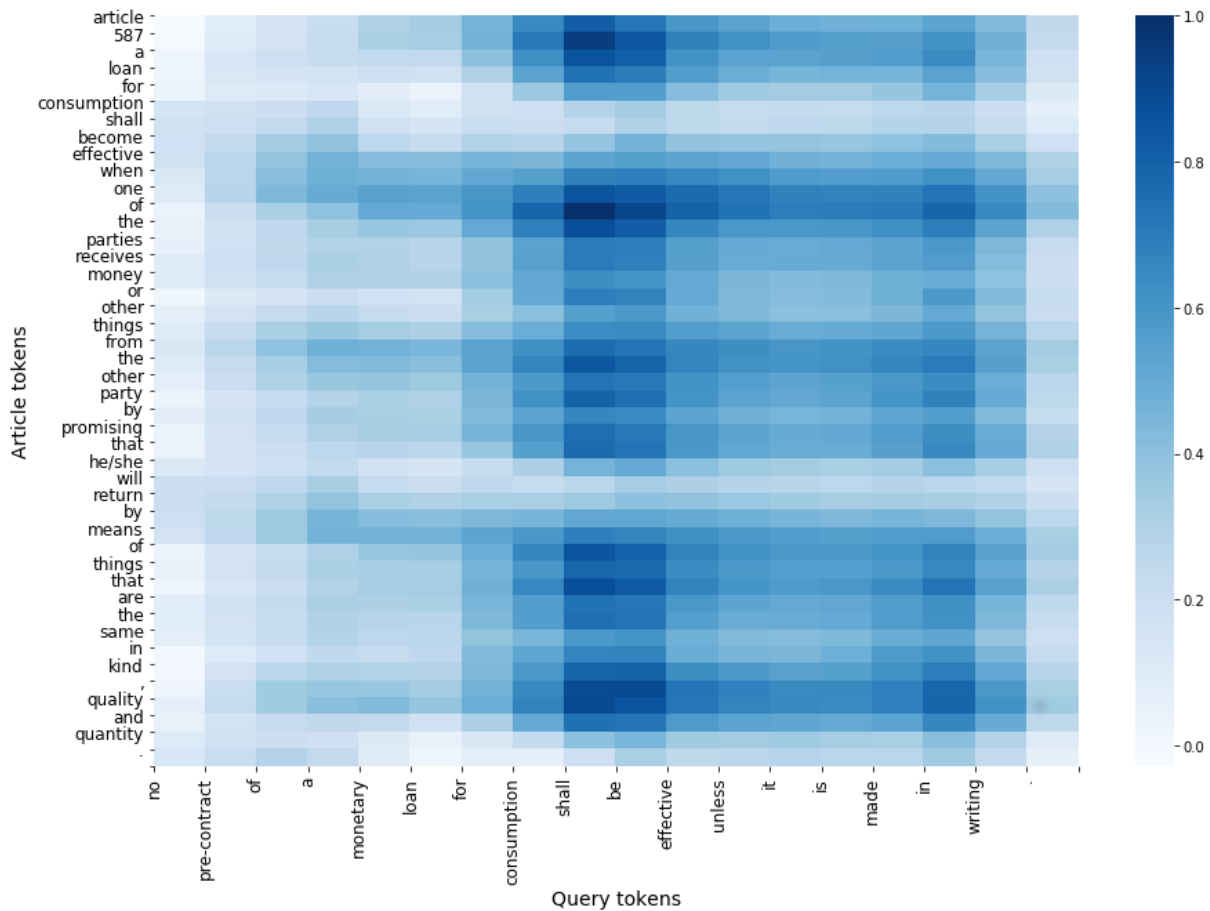


Figure 3: Neural attention example for task H30-23-A. Our model incorrectly predicted a negative entailment relationship to be positive.

Since our model works with neural attention, we can visualize the query-document interaction in our model for selected cases. For instance, in task H30-23-A, our model predicted the positive class, although the correct answer is a negative entailment.

Figure 3 illustrates the neuron activation of one of our stacked encoder models. Note that there is almost no neuron activation between the query term “pre-contract” and the article tokens. We suppose that the Glove embeddings do not provide feature information for the legal use of this specific hyphenated term. Interestingly, there is a high neuron activation between the query sequence “shall be effective” and “when one of the parties receives”, as well as the article number within the law document. The sequence “that are the same in kind” also exhibits higher activation values for that query sequence. Therefore, we attribute this mistake to the out-of-vocabulary-problem, which is often encountered in domain-specific texts with scarce resources for context modeling with embeddings. While own Word2Vec embeddings were employed for the information retrieval task, we selected Glove embeddings because of the presumably higher likelihood of correctly learned semantic feature representations for negation words and terms such as “must” and “may”.

The second example in Figure 4 shows a correctly classified positive entailment relationship. Elevated activation scores are found in the query terms “no successor of the primary obligor may” in conjunction with the article term sequence “no primary obligor, guarantor or successor” and “may make a claim for the extinction”. In this case, it can also be observed that the article number has not been attended as strong as in the previous example.

The previous two examples show that our stacked encoder has the ability to attend meaningful inputs, given that the input feature representation is expressive enough to capture the domain-specific context of a word. Due to space limitations, examples with positive predictions were chosen because both encoders are highly biased towards the negative class. In the few cases where they predict a positive entailment, they appear to have a local competence for correctly identifying vital features. Therefore, we expect the performance to increase when word embeddings are generated from a large domain-specific corpus.

6 CONCLUSION

In this work, we present a law retrieval approach and an entailment classifier, developed during the COLIEE ’19 competition. For

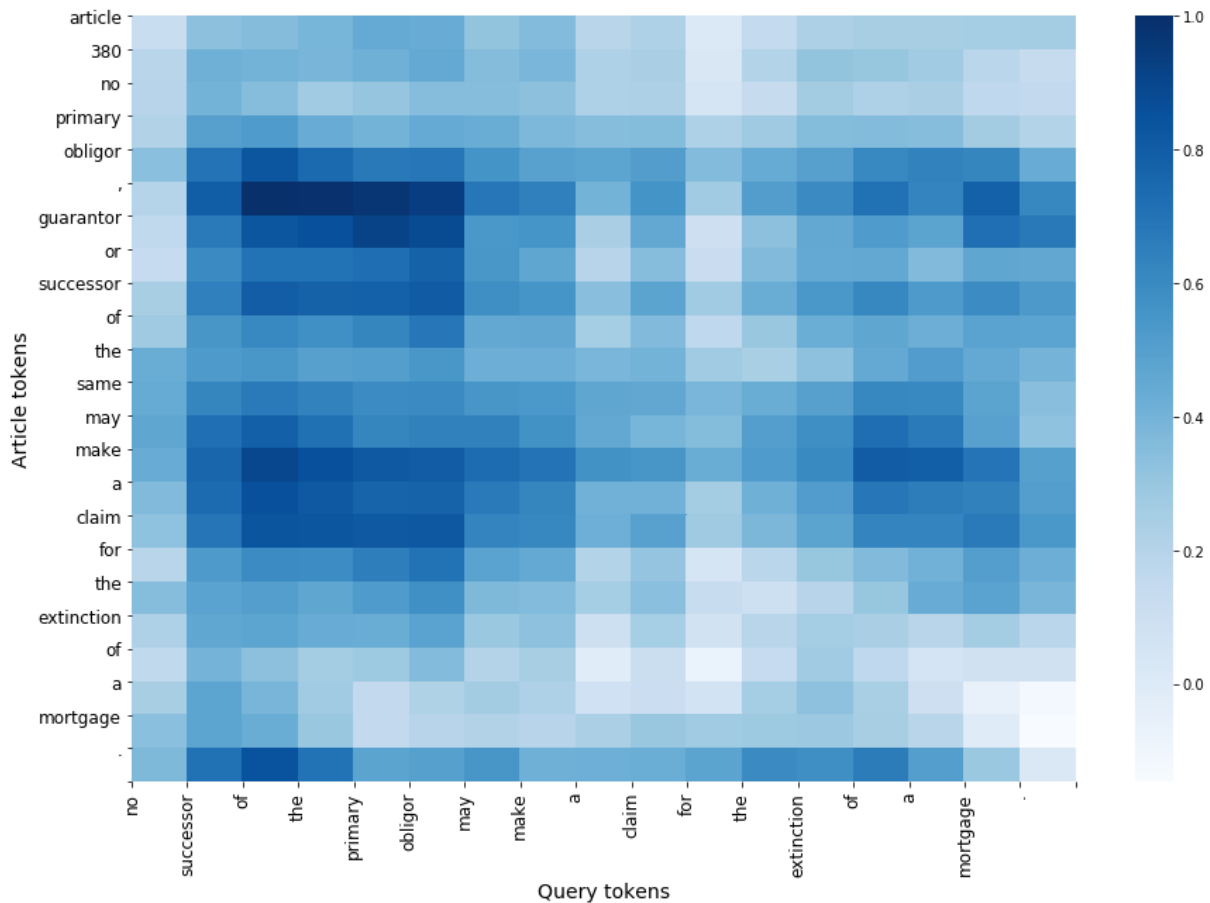


Figure 4: Neural attention example for task H30-13-I. Our model correctly predicted a positive entailment relationship.

the law retrieval task, we combined BM25 scoring with the Word Centroid Distance for word embeddings, and then applied similarity thresholding for a variable number of retrieved documents per query. Our results outperform some competitor approaches, but we expect further improvements by using word embeddings trained on a larger legal corpus. Regarding the textual entailment task, we outperformed the baseline and thereby have shown a benefit of using a stacked encoder architecture with additional features and an attention layer. Future work for the entailment task can be done on additional feature extraction, for instance semantic role labels. Since our ensemble has been selected manually, it can be beneficial to use automated pruning approaches which are designed for imbalanced data to overcome classifier bias and exploit complementary local competences [14]. Another research field is the use of different learning architectures and other word embeddings suitable for the legal domain.

REFERENCES

- [1] Piyush Arora, Murhaf Hossari, Alfredo Maldonado, Gareth J.F. Conran, Clare and Jones, Johannes Paulus, Alexander andKlostermann, and Christian Dirschl. 2018. Challenges in the Development of Effective Systems for Professional Legal Search. In *ProfS/KG4IR/Data: Search@ SIGIR*. CEUR-WS.org, Ann Arbor, Michigan, USA, 29–34.
- [2] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473* (2014).
- [3] Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Lisbon, Portugal, 632–642. <https://doi.org/10.18653/v1/D15-1075>
- [4] Danilo S. Carvalho, Minh-Tien Nguyen, Chien-Xuan Tran, and Minh-Le Nguyen. 2017. Lexical-Morphological Modeling for Legal Text Analysis. In *New Frontiers in Artificial Intelligence*, Mihoko Otake, Setsuya Kurahashi, Yuiko Ota, Ken Satoh, and Daisuke Bekki (Eds.). Springer International Publishing, Cham, 295–311.
- [5] Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. Reading Wikipedia to Answer Open-Domain Questions. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Vancouver, Canada, 1870–1879. <https://doi.org/10.18653/v1/P17-1171>
- [6] Phong-Khac Do, Huy-Tien Nguyen, Chien-Xuan Tran, Minh-Tien Nguyen, and Minh-Le Nguyen. 2017. Legal question answering using ranking SVM and deep convolutional neural network. *arXiv preprint arXiv:1703.05320* (2017).
- [7] Lukas Galke, Ahmed Saleh, and Ansgar Scherp. 2017. Word Embeddings for Practical Information Retrieval. *INFORMATIK 2017* (2017), 2155–2167.
- [8] Oren Glickman and Ido Dagan. 2005. Web based probabilistic textual entailment. In *In Proceedings of the 1st Pascal Challenge Workshop*. 33–36.
- [9] Quang-Thuy Ha, Thi-Oanh Ha, Thi-Dung Nguyen, and Thuy-Linh Nguyen Thi. 2012. Refining the Judgment Threshold to Improve Recognizing Textual Entailment Using Similarity. In *Computational Collective Intelligence. Technologies and Applications*, Ngoc-Thanh Nguyen, Kiem Hoang, and Piotr Jedrzejowicz (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 335–344.

- [10] Aminul Islam and Diana Inkpen. 2008. Semantic text similarity using corpus-based word similarity and string similarity. *ACM Transactions on Knowledge Discovery from Data (TKDD)* 2, 2 (2008), 10.
- [11] Mi-Young Kim, Randy Goebel, Yoshinobu Kano, and Ken Satoh. 2016. COLIEE-2016: evaluation of the competition on legal information extraction and entailment. In *International Workshop on Juris-informatics (JURISIN 2016)*.
- [12] Mi-Young Kim, Randy Goebel, and Ken Satoh. 2015. COLIEE-2015: evaluation of legal question answering. In *Ninth International Workshop on Juris-informatics (JURISIN 2015)*.
- [13] Mi-Young Kim, Yao Lu, Juliano Rabelo, and Randy Goebel. 2018. COLIEE-2018: Evaluation of the Competition on Case Law Information Extraction and Entailment. [online]. https://sites.ualberta.ca/~rabelo/COLIEE2019/COLIEE2018_CL_summary.pdf
- [14] Bartosz Krawczyk and Michał Woźniak. 2018. Leveraging Ensemble Pruning for Imbalanced Data Classification. In *2018 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*. IEEE, 439–444.
- [15] Matt J. Kusner, Yu Sun, Nicholas I. Kolkin, and Kilian Q. Weinberger. 2015. From Word Embeddings to Document Distances. In *Proceedings of the 32Nd International Conference on International Conference on Machine Learning - Volume 37 (ICML '15)*. JMLR.org, 957–966. <http://dl.acm.org/citation.cfm?id=3045118.3045221>
- [16] Yang Liu, Chengjie Sun, Lei Lin, and Xiaolong Wang. 2016. Learning natural language inference using bidirectional LSTM model and inner-attention. *arXiv preprint arXiv:1605.09090* (2016).
- [17] Arpan Mandal, Kripabandhu Ghosh, Arnab Bhattacharya, Arindam Pal, and Saptarshi Ghosh. 2017. Overview of the FIRE 2017 IRLeD Track: Information Retrieval from Legal Documents. In *FIRE (Working Notes)*, Prasenjit Majumder, Mandar Mitra, Parth Mehta Mehta, and Jainisha Sankhavara (Eds.). CEUR-WS.org, Bangalore, India, 63–68.
- [18] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. 2008. *Introduction to Information Retrieval*. Cambridge University Press, New York, NY, USA.
- [19] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Distributed Representations of Words and Phrases and Their Compositionality. In *Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 2 (NIPS'13)*. Curran Associates Inc., USA, 3111–3119. <http://dl.acm.org/citation.cfm?id=2999792.2999959>
- [20] Rohan Nanda, Adebayo Kolawole John, Luigi Di Caro, Guido Boella, and Livio Robaldo. 2017. Legal Information Retrieval Using Topic Clustering and Neural Networks. In *COLIEE 2017. 4th Competition on Legal Information Extraction and Entailment (EPIC Series in Computing)*, Ken Satoh, Mi-Young Kim, Yoshinobu Kano, Randy Goebel, and Tiago Oliveira (Eds.), Vol. 47. EasyChair, 68–78. <https://doi.org/10.29007/psgx>
- [21] Yixin Nie and Mohit Bansal. 2017. Shortcut-Stacked Sentence Encoders for Multi-Domain Inference. In *Proceedings of the 2nd Workshop on Evaluating Vector Space Representations for NLP*. Association for Computational Linguistics, Copenhagen, Denmark, 41–45. <https://doi.org/10.18653/v1/W17-5308>
- [22] Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. Glove: Global Vectors for Word Representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Doha, Qatar, 1532–1543. <https://doi.org/10.3115/v1/D14-1162>
- [23] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ Questions for Machine Comprehension of Text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Austin, Texas, 2383–2392. <https://doi.org/10.18653/v1/D16-1264>
- [24] Siva Reddy, Danqi Chen, and Christopher D Manning. 2018. Coqa: A conversational question answering challenge. *arXiv preprint arXiv:1808.07042* (2018).
- [25] Stephen Robertson and Hugo Zaragoza. 2009. The Probabilistic Relevance Framework: BM25 and Beyond. *Information Retrieval* 3, 4 (2009), 333–389.
- [26] Tim Rocktäschel, Edward Grefenstette, Karl Moritz Hermann, Tomáš Kočiský, and Phil Blunsom. 2015. Reasoning about entailment with neural attention. *arXiv preprint arXiv:1509.06664* (2015).
- [27] Niall Rooney, Hui Wang, and Philip S Taylor. 2014. An investigation into the application of ensemble learning for entailment classification. *Information Processing & Management* 50, 1 (2014), 87–103.
- [28] Johannes Schmidt-Hieber. 2017. Nonparametric regression using deep neural networks with ReLU activation function. *arXiv preprint arXiv:1708.06633* (2017).
- [29] Benno Stein, Sven Meyer zu Eissen, and Martin Potthast. 2007. Strategies for Retrieving Plagiarized Documents. In *Proceedings of the 30th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '07)*. ACM, New York, NY, USA, 825–826. <https://doi.org/10.1145/1277741.1277928>
- [30] Eylon Stroh and Priyank Mathur. 2016. Question Answering Using Deep Learning. [online]. <https://cs224d.stanford.edu/reports/StrohMathur.pdf>
- [31] Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. Sequence to Sequence Learning with Neural Networks. In *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2 (NIPS'14)*. MIT Press, Cambridge, MA, USA, 3104–3112. <http://dl.acm.org/citation.cfm?id=2969033.2969173>
- [32] Liuyang Tian, Hui Ning, Leilei Kong, Zhongyuan Han, Ruiming Xiao, and Hao-liang Qi. 2017. HLJIT2017@ IRLeD-FIRE2017: Information Retrieval From Legal Documents. In *FIRE (Working Notes)*, Prasenjit Majumder, Mandar Mitra, Parth Mehta Mehta, and Jainisha Sankhavara (Eds.). CEUR-WS.org, Bangalore, India, 82–85.
- [33] Dmitry Yarotsky. 2018. Optimal approximation of continuous functions by very deep ReLU networks. *arXiv preprint arXiv:1802.03620* (2018).
- [34] Masaharu Yoshioka, Yoshinobu Kano, Naoki Kiyota, and Ken Satoh. 2018. Overview of Japanese Statute Law Retrieval and Entailment Task at COLIEE-2018. [online]. https://sites.ualberta.ca/~rabelo/COLIEE2019/COLIEE2018_SL_summary.pdf