



School of Information Sciences and Technology
Department of Informatics
Athens, Greece

Bachelor Thesis
in
Computer Science

Data Augmentation for Legal Document Classification

Platon Karageorgis

Supervisors: Prof. Ion Androutsopoulos

Department of Informatics
Athens University of Economics and Business

Dr. Dimitris Pappas

Department of Informatics
Athens University of Economics and Business

September 2023

Platon Karageorgis

Data Augmentation for Legal Document Classification

September 2023

Supervisors: Prof. Ion Androutsopoulos, Dr. Dimitris Pappas

Athens University of Economics and Business

School of Information Sciences and Technology

Department of Informatics

Natural Language Processing Group

Athens, Greece

Abstract

Data augmentation is a technique for generating artificial instances of data. It is widely known in Computer Vision and it has been used for years in Natural Language Processing (NLP). Despite its widespread use, its application in the legal domain remains relatively unexplored. Legal texts, characterized by intricate language and lengthy documents, present unique challenges that demand specialized solutions. We implement a variety of augmentation methods aiming to demonstrate that data augmentation can enhance significantly the performance of language models in legal document classification. LEGAL-BERT, a specialized language model designed for legal text, serves as the foundation for our experiments. We demonstrate that the legal domain presents complexities beyond initial observation, and we attempt various approaches to overcome the barriers that are raised during our experiments.

Περίληψη

Η επαύξηση δεδομένων είναι μία τεχνική που στοχεύει στη γέννηση νέων, τεχνητών δεδομένων. Είναι ευρέως γνωστή τεχνική σε εκείνους που ανήκουν στο χώρο της υπολογιστικής όρασης και χρησιμοποιείται χρόνια στο χώρο της επεξεργασίας φυσικής γλώσσας. Παρά την πληθώρα των εφαρμογών της στον τομέα αυτό σε γενικότερο πλαίσιο, είναι γεγονός ότι η συνεισφορά της σε νομικά κείμενα δεν έχει διερευνηθεί σε ικανοποιητικό βαθμό. Τα νομικά κείμενα χαρακτηρίζονται από την εκτενή χρήση πολύπλοκου λεξιλογίου αλλά και από τη μεγάλη έκτασή τους. Έτσι, απαιτούν ειδική αντιμετώπιση, όσον αφορά την επεξεργασία φυσικής γλώσσας αλλά και κατ'επέκταση την επαύξηση δεδομένων. Υλοποιούμε μία σειρά από τεχνικές επαύξησης δεδομένων, προσπαθώντας να δείξουμε τη θετική συνεισφορά που μπορεί να έχει η επαύξηση δεδομένων στα γλωσσικά μοντέλα που πραγματοποιούν ταξινόμηση νομικών κειμένων. Παράλληλα, επιδεικνύουμε και αποδεικνύουμε τη βαθύτερη πολυπλοκότητα που χαρακτηρίζει τα νομικά κείμενα και ελυσόμαστε με μία σειρά από τεχνικές προκειμένου να αντιμετωπίσουμε τα εμπόδια που προκύπτουν κατά τη διάρκεια των πειραμάτων μας.

Acknowledgements

I would like to thank Professor Ion Androutsopoulos for giving me the opportunity to work in an engaging field like NLP alongside a team of professionals. I am humbled by his trust and patient guidance throughout all these months and I hope he will remain an inspiration for the next NLP enthusiasts to come. Moreover, I would like to express my heartfelt thanks to Dr. Dimitris Pappas, who was there whenever I needed him to answer my questions and provide me with guidance. His contribution was instrumental to my growth as an aspiring researcher. I am also deeply appreciative of Dr. Prodromos Malakasiotis for his state-of-the-art insights, which elevated the caliber of this thesis. Last but not least, as this thesis marks the end of my Bachelor's Degree at Athens University of Economics and Business, I would like to thank the entire faculty and staff for the support they provided over the past five years.

Contents

Abstract	iv
Acknowledgements	vii
1 Introduction	1
1.1 Motivation and Problem Statement	2
1.2 Thesis Structure	4
2 Data Augmentation Methods	5
2.1 Data Augmentation in Biomedical Domain	5
2.1.1 Models	5
2.1.2 Word2vec	6
2.1.3 Masked-LM	6
2.1.4 Back-Translation	7
2.1.5 Additional Augmentation Methods	7
2.2 Data Augmentation In Legal Domain	8
2.2.1 Our Work	8
2.2.1.1 Word2vec	8
2.2.1.2 Masked-LM	10
2.2.1.3 Back-Translation	12
2.2.2 Baselines	13
2.2.2.1 Naive Oversampling Baseline	13
2.2.2.2 Masked-LM without Substitution	13
2.2.3 Data Augmentation in Legal Texts	13
2.2.3.1 Random Deletion/Insertion/Scrambling/Swapping	14
2.2.3.2 Data Augmentation using Semantic Relations	15
2.2.4 Additional Data Augmentation Methods in NLP	16
2.2.4.1 Synonym Replacement	16
2.2.4.2 Data Augmentation using LLMs	17
3 Experiments	19
3.1 Experimental Setup	19
3.1.1 Hardware Setup	19
3.1.2 SCOTUS	20
3.1.3 LEGAL-BERT	22

3.1.4	Hierarchical Model	23
3.1.5	Evaluation Metrics	25
3.2	Experimental Results	27
3.2.1	SCOTUS Results	27
3.2.1.1	Basic & Smart Augmentation	27
3.2.1.2	Visualization & Analysis	33
4	Conclusions and Future Work	39
4.1	Conclusions	39
4.2	Future Work	40
	Bibliography	41
	List of Figures	44
	List of Tables	46

Introduction

The rise of Artificial Intelligence and mainly Natural Language Processing (NLP), especially during the last five years, has been outstanding. A milestone for NLP was the shift from Convolutional Neural Networks and Recurrent Neural Networks to Transformers. This led to the development of BERT [Dev+18] which is a type of pre-trained Transformer and remains to this day a very powerful tool. The latest update is the ascension of Large Language Models (LLMs), with the lead star being Chat-GPT, whose capabilities have traveled to households of every cornerstone of the planet. At the time that this thesis is written, LLMs are admittedly the hot topic that every AI researcher has to mention in any NLP-related task.

Despite the surprisingly good results of LLMs across a wide spectrum of applications, in cases where smaller non-generative models like BERT are viable, especially for specific tasks such as classification, they just might be the runner-up. The challenge that arises though when BERT is fine-tuned, is a common barrier in the Deep Learning era. The amount of annotated data needed to fine-tune BERT up to the point where the evaluation measures (usually F1-score) are optimized, most of the time surpassing the realistic number of data available in a real-world task. Legal documents are a great example of this, as they have long but few texts on many occasions. Data augmentation is a technique that increases the amount of available annotated data, by altering the existing ones or creating entirely new data rows from scratch. It is widely known in Computer Vision as it is more straightforward to implement when dealing with images. On the other hand, it is not so simple to apply it in NLP, as text documents are not as easily augmented.

For example, consider the sentence "The court decided to postpone the hearing for next week". The most intuitive manner of augmentation would be to replace several words with their synonyms to generate a different data row that also maintains the original meaning. For instance, the above sentence could be transformed to "The court chose to put off the hearing for next week". In this case, the augmentation process was effectively implemented, but is it really easy to automate this synonym replacement process using a language model? The word "decided" has multiple synonyms in general, but in this specific context, a potentially big subset of them will not be a fitting choice. This adds a whole new level of complexity in text augmentation that sources from the perplexity of natural languages. Ultimately, it can be observed that even when using the most intuitive augmentation method, text document augmentation is more complex than image

augmentation. The topic of this thesis covers the effect of a selection of data augmentation techniques in legal document classification.

1.1 Motivation and Problem Statement

Data augmentation is likely familiar to those who specialize in Computer Vision. Image processing is a domain that people have been working on far before Transformers or Neural Networks were developed. And that is all it takes to augment data in Computer Vision, as the data rows are not texts, but images. A mere rotation, a slight alteration on the pixel settings, or an even more complicated method like image compression using Singular Value Decomposition (SVD) truncation as shown in [Figure 1.1](#), could be enough to generate artificial data. Conversely, the application methods of Data augmentation in Natural Language Processing vary and have quite different backgrounds due to the natural difference between text and images. The prime inspiration of this thesis is Pappas et al. [PMA22]. The paper considered several data augmentation methods in the context of biomedical question-answering, which included an artificial cloze-style machine reading comprehension dataset, back-translation, information retrieval, word substitution using word2vec and masked-lm, question generation, and context addition. The results of the experiments were very positive and showed that data augmentation can be a key contributing factor in improving the performance of a model. More details regarding the work on this paper will be provided in [Section 2](#) that will cover related work on various data augmentation techniques. In this thesis, we will inherit a subset of these augmentation methods and apply them in legal text classification. The augmentation techniques that will be inherited from the mentioned paper are both word substitution implementations, meaning word2vec and masked-lm, and finally back-translation. These methods will be extensively analyzed in subsequent sections.

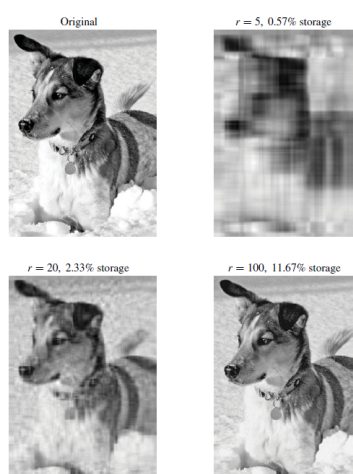


Fig. 1.1: Image compression of Mordecai the snow dog, truncating the SVD at various ranks r and therefore generating artificial image data. Figure copied from the book *Data-Driven Science and Engineering 2nd Edition*, by Steven L. Brunton and J. Nathan Kutz

As we mentioned in the introductory paragraph, the previously stated techniques will be implemented in legal document classification. The dataset will be drawn from LexGLUE [Cha+22]. LexGLUE is a benchmark collection of datasets created to evaluate models in legal tasks. The files and paperwork needed for legal tasks have always been a factor that makes law a challenging sector to contribute to. AI and NLP aim to take this a step further by securing automation in tasks that have been troubling legal practitioners. LexGLUE makes a step in that direction by designing a benchmark that will benefit researchers in Natural Language Understanding (NLU) tasks by improving model performance and evolving their techniques through this process. In this thesis, we use SCOTUS, a part of the LexGLUE dataset for multi-class classification. The selection was based on the amount of available data, as we expect datasets with less available data to get an additional performance boost using more synthetic (augmented) data. More information regarding the work on LexGLUE will be once again provided in [Section 3](#) that describes our experiments but also gives a quick background for the adopted datasets.

Moreover, the dataset described above will be augmented using various techniques that will be extensively analyzed in [Section 2](#). Subsequently, original and artificial data will be merged, to train the chosen language model. The model that will be used in all of the experiments in this thesis will be LEGAL-BERT, constructed by Chalkidis et al. [Cha+20]. The work of Chalkidis et al. acknowledges the fantastic results that BERT achieves in NLP but also underlines that it fails to produce equally good results in specific scientific areas. Therefore, the authors attempt to manufacture a version that will be focused solely on the legal domain. Taking into account the dataset that will be used for the experiments, LEGAL-BERT seems like an excellent choice for our task. It was pre-trained on legal data and, therefore, is expected to perform exceptionally well in the legal domain. In fact, it is not only one of the models that LexGLUE used to assess the generated benchmark dataset, but it is also the best-performing one in the majority of the metrics. To conclude, we aim to improve upon the LEGAL-BERT results reported by Chalkidis et al. [Cha+20] by applying data augmentation.

Having analyzed briefly the necessary background, we can now state the problem that this thesis will be devoted to. Data augmentation is a technique that generates artificial data and it has been proved in past work that it can be quite beneficial in improving model results. The biomedical domain is a great example of this, as Pappas et al. [PMA22] proved through their positive results. Using the results of LexGLUE, a legal benchmark dataset, we apply a variety of data augmentation methods to prove that data augmentation can be beneficial in the legal domain as well. The model used to perform the experiments is LEGAL-BERT, the best model overall in LexGLUE.

1.2 Thesis Structure

The remainder of this thesis is structured as follows:

Chapter 2: Data Augmentation Methods

The next chapter of this thesis will cover at first the data augmentation methods applied in the biomedical domain by Pappas et al. [PMA22] which is the paper that we inherit our methods. Then, there will be a discussion on the techniques put into effect in this thesis, along with other implementations in the legal domain. Lastly, the final section will concern a few augmentation methods outside the spectrum of specialized scientific areas, like the biomedical and the legal ones.

Chapter 3: Experiments

The third chapter will first present the background of the dataset and the model used for the experiments, as well as the metrics and the baselines that were adopted to support them. Subsequently, there will be a detailed presentation of the experiments that were carried out, along with their corresponding results. These experiments aim to compare the effects of our techniques with the benchmark results produced by the original, non-augmented dataset of LexGLUE. Moreover, there will be a discussion on the quality of the results.

Chapter 4: Conclusions

The fourth and final chapter will primarily sum up the goals of this thesis and comment on the results of the experiments. Finally, there will be suggestions for future work on data augmentation in the legal domain, but also in NLP in general.

Data Augmentation Methods

This section is separated into three parts. The first part explores the seven augmentation methods implemented in the biomedical domain proposed by Pappas et al. [PMA22]. The second part tackles data augmentation in the legal domain and is further segmented into two components. The initial one covers our work and the latter discusses similar augmentation methods in a legal context, implemented by the research community. Finally, the third section introduces a couple of extra augmentation methods generally in NLP, not necessarily concerning the legal or biomedical domain.

2.1 Data Augmentation in Biomedical Domain

In 2022, Pappas et al. [PMA22] investigated the impact of seven data augmentation techniques in the context of factoid question answering, with a specific emphasis on the biomedical domain. A selection of models in the BERT family was recruited for this cause, and the employed data were drawn from the BIOASQ challenge [Tsa+15]. It's pertinent to note that the data have a specific shape, each data row is a question-snippet-answer triple. As mentioned in the introduction we have inherited a subset of these augmentation techniques and we will assess their effect on the legal domain. Before continuing with this assessment though, we will focus on the work done in a biomedical context, and we will describe the aspects of each augmentation technique, as some of these are implemented in this thesis.

2.1.1 Models

Before delving into the analysis of augmentation methods, let's introduce the models that will be used throughout the paper. The authors use DISTILBERT [San+20], BIOBERT [Lee+19] and ALBERT [Lan+20] which are BERT-like pre-trained models. DISTILBERT is a lighter and faster version of the basic version of BERT, bert-base-uncased. BIOBERT is a version of BERT pre-trained on biomedical documents, and ALBERT is like DISTILBERT, a lighter version of BERT as it attempts to limit BERT's computational demands. They are fine-tuned on SQUAD [Raj+16] or SQUAD-v2 [RJL18] which are Stanford's question-answering datasets, and ALBERT is further fine-tuned on BIOASQ data aiming to be a strong baseline. There are further modifications to the models that will not be mentioned, as we aim to focus on the augmentation techniques.

2.1.2 Word2vec

To begin with, the first augmentation method is implemented using word2vec, a well-known tool for those that specialize in NLP, initially proposed by Mikolov et al. [Mik+13]. Word2vec represents words in a vector space where each word is depicted by a unique vector known as a word embedding.

The research team of Pappas et al. [PMA22] analyzed question-snippet-answer training instances and began by examining the word tokens in the snippet, excluding stop-words. For each token w_i within the snippet, ranging from 1 to n , the authors chose up to $k_i \leq K$ (where K is the upper bound of the total words to be retrieved) most similar words w_j from the vocabulary based on cosine similarity between word embeddings \vec{w}_i and \vec{w}_j (with $\cos(\vec{w}_i, \vec{w}_j) \geq C$, where C is the similarity threshold). This process resulted in generating $(k_1+1)(k_2+1)\dots(k_n+1)-1$ artificial training instances, with at least one token changed in each instance. It's worth noting that the total number of neighboring words k_i (which is set at 10) for a given word may not always reach the specified value. This is due to the authors' choice of a stringent 0.95 similarity threshold, which results in a smaller pool of words exceeding this threshold compared to k_i .

Furthermore, the authors sampled training instances ranging from 10,000 to 100,000 and integrated them into the training dataset, contributing to a substantial increase in the F1-score by up to six percentage points. Finally, the decision to set the values for the parameters K and C (similarity threshold) to 10 and 0.95, respectively, was informed by preliminary experiments conducted on development data.

2.1.3 Masked-LM

The second augmentation method is masked-lm. Like word2vec, masked-lm is a word substitution technique but operates using a different scheme. It takes advantage of the fact that BERT-like models are trained for Masked Language Modelling. These models are exposed to millions of sentences with hidden words and are tasked with identifying the missing token. The authors employ BIOLM [Lew+20], specifically a ROBERTA-LARGE model pre-trained on PUBMED Baseline Repository¹, PMC, and MIMIC-III [ZPT18], enriched with a new vocabulary sourced from PUBMED.

The difference in this approach in comparison with word2vec is that each potential replacement word w_j for an original word w_i in the snippet must meet the condition $p(w_j) \geq P$ (instead of $\cos(w_i, w_j) \geq C$ enforced in word2vec). The term $p(w_j)$ represents the probability assigned to w_j by the pre-trained model. Additionally, they arrange the candidate

¹See <https://pubmed.ncbi.nlm.nih.gov>

replacements w_j for each w_i based on their $p(w_j)$ values. They establish $P=0.95$, grounded in preliminary experiments conducted on development data.

It's noteworthy that in their experiments, models trained using BIOLM for data augmentation performed similarly to models trained using word2vec for data augmentation, but with an important difference. For BIOLM, the optimal performance is achieved when incorporating 50k artificial examples, in contrast to 10k for word2vec. The authors hypothesize that this happens due to BIOLM suggesting words that align closely with the specific context of the word being replaced, whereas word2vec takes a more direct approach by searching candidate replacements for each original word.

2.1.4 Back-Translation

The third augmentation technique is back-translation, a technique widely applied in NLP tasks [SKF21], [Fen+21]. This method involves translating training examples from a source language to a pivot language and then back to the source language, thus generating paraphrased versions of the original text. Initially, the authors of Pappas et al. [PMA22] utilized French as the pivot language and later added Spanish and German.

For each question-snippet-answer training triplet within BIOASQ, there are two additional triplets created by performing back translation on either the question or the snippet. If any of the new triplets happen to match the original, they are discarded. This process yields 4,901 new training examples when pivoting solely to French, and 15,593 when extending the pivot languages to include Spanish and German.

Adding the generated data into the BIOASQ training dataset results in a nearly 2% enhancement in the development data F1-score, regardless of using one or three pivot languages. Also, the authors attempted to run the model only with artificial data and the results were almost identical with the original ones, proving therefore the quality of artificial data. Finally, the authors mention that simpler methods like word2vec, can yield more substantial gains with fewer artificial training instances, and perhaps are a better choice regarding augmentation methods.

2.1.5 Additional Augmentation Methods

Up to this point, we have discussed techniques that will be further explored in this thesis, to evaluate their performance in a legal context. The following four methods will not be the focus of our work. Firstly, the authors implement an artificial cloze-style Machine Reading Comprehension (MRC) augmentation method that leverages biomedical articles to generate question-snippet-answer triples, significantly expanding the dataset. Secondly, the authors

use information retrieval systems to extract relevant content from sources like PUBMED. The third technique employs T5 [Raf+23] for question generation, creating new question-snippet-answer triples from BIOASQ training data, with modifications to highlighted areas. Lastly, context addition adds neighboring sentences to snippets, enhancing context awareness. It uses variables to determine the extent of sentence addition, with a focus on maintaining character limits.

All of these approaches achieved results surpassing the established baseline and were deemed successful by the authors of Pappas et al. [PMA22]. However, it's worth noting that the information retrieval method yielded less data than initially expected, which influenced our decision to exclude it from our thesis. Question generation, too, does not align with the nature of our task, which primarily focuses on classification rather than question generation. Additionally, the artificial cloze-style Machine Reading Comprehension (MRC) augmentation and context-aware methods are not suitable for our dataset, which is centered around LexGLUE benchmark datasets, and these techniques typically involve snippets extracted from articles.

2.2 Data Augmentation In Legal Domain

Data augmentation in legal texts is not widely recognized as a prominent domain within the field of Natural Language Processing (NLP). Nevertheless, in recent years researchers have demonstrated an increased interest in the field as Csányi et al. [CO21] described. In this section, we refer to relevant work in the field, but before that, we focus on the methods implemented in this thesis.

2.2.1 Our Work

Our work consists of three main methods that perform data augmentation. They all have been inherited from Pappas et al. [PMA22] and we shortly discussed them above in a biomedical context. In this part of the thesis, we will analyze word2vec, masked-lm, and back-translation once more, but from our scope this time, adding more details and examples.

2.2.1.1 Word2vec

Text augmentation using word2vec is essentially an advanced word substitution. Each word in the text is examined using pre-trained legal embeddings called Law2Vec [CK19]. Law2Vec's extensive corpus comprises over 123,000 documents with 492 million tokens. Texts were preprocessed by removing non-UTF8 characters and tokenizing using NLTK.

Additionally, all words were converted to lowercase, and numbers were replaced with 'D' [CA17] to maximize matches during the embedding search.

When searching for word embeddings, there are two outcomes. If a corresponding embedding is not found, the word remains unchanged in the augmented text. However, if the embedding is found, the method retrieves its vector representation and searches for its n most similar embeddings using cosine similarity. The formula is the same as in the biomedical implementation, $\cos(\vec{w}_i, \vec{w}_j) \geq C$ (where C is the similarity threshold). These embeddings constitute potential candidates to replace the initial word, forming a 'bag of words' for each word in the text. It's important to note that the initial word is also included in its bag, which means it may still appear in the augmented text, even in the case where other candidates different than the original word are available, as seen in [Figure 2.1](#).

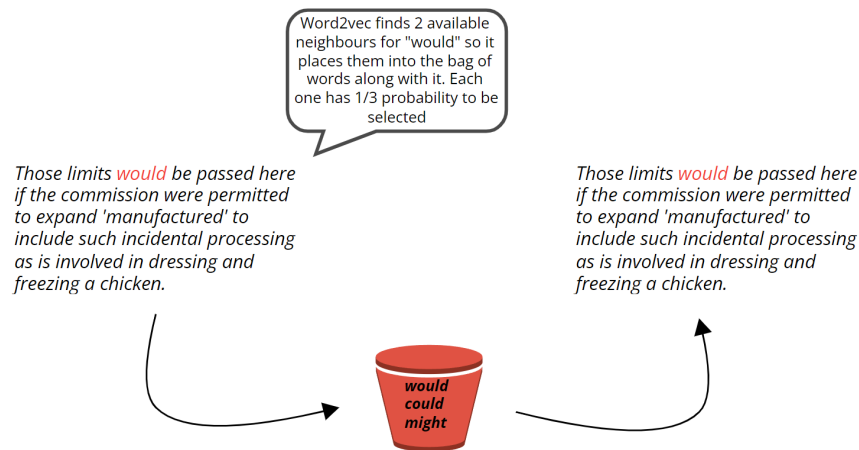


Fig. 2.1: Word2vec augmentation example in a chunk of our SCOTUS legal texts.

As we mentioned in the previous description of word2vec in a biomedical context, it's not certain whether the bucket will contain in total n neighboring words for a given word. This occurs because for the method to be precise, the resulting words must truly be "close" to the target word to preserve the overall sentence meaning. This is achieved by setting a threshold that will set how different the words are allowed to be compared to the initial one. We set the threshold C to be 0.95 after tuning it using the development dataset. This means we allow the word to deviate from the original by up to 5%. For a random word, the bag of words created via this process can have n words, but it can also have less or even none. In practice, due to the high threshold, the latter case will likely be more common.

Finally, after processing the entire text and creating the bags of words for each word, the augmentation process begins. For each word in the text, a replacement word is ran-

domly selected from its corresponding bag. This process can be repeated indefinitely, potentially generating a large number of text variations. However, it's essential to note that after a certain point, the generated texts may become identical to those previously created. The maximum number of possible text variations can vary and depends on several factors, including the similarity threshold, text length, and randomness. Also, some texts may contain words more likely to have suitable replacement candidates due to their inherent characteristics.

2.2.1.2 Masked-LM

The implementation of masking using a language model is relatively straightforward, with one exception. Due to our documents being way too long for the model to process (some of the biggest documents have over 180,000 characters), we have to make repeated segmentations before feeding the text to the model. To enhance the effectiveness of our code, we utilize a tailored text-breaking technique, which will be described in detail below.

Our approach to text segmentation is as follows. Knowing that LEGAL-BERT has a token limit of 512, we set the maximum number of BPEs, which represent sub-word units, to be 500. In a few words, Byte-Pair Encodings (BPEs) is a sub-word tokenization technique used to break down words into smaller units for natural language processing tasks, such as machine translation and text generation. They were initially proposed by Sennrich et al. [SHB16]. The choice of setting the maximum number of BPEs to 500, accounts for potential BPE tokenization that may split a token into two, causing the total token count to exceed 512. Therefore, we set the threshold slightly lower to avoid exceeding the limit. As we process the document, the reading process pauses when the token limit, based on BPEs, is reached. We then backtrack through the text to identify the last paragraph change before the point where the limit was exceeded. If a paragraph change is found, we split the text at that point, and this chunk becomes a paragraph. If no paragraph change is detected, we repeat the process by looking for full stops. In rare cases, such as a document with no full stops or paragraph changes within a 500-token span, we break the text at the last line change. This process is repeated for the next chunk, and we meticulously track the type of segmentation that occurred between each pair of chunks. The goal is to assemble them at the end of the iteration, once the entire document has been processed. This careful process ensures that the final augmented data row maintains up to one difference in each chunk from the original text, which corresponds to the predicted word. This approach guarantees that the rest of the text remains identical to the original.

Once we have a valid chunk in terms of length, the subsequent part is masked-lm prediction. A token contained in the chunk is randomly chosen and replaced by a [MASK] token. Then, the chunk is tokenized using LEGAL-BERT and afterward, we ask the model to

predict the word behind [MASK]. LEGAL-BERT replaces [MASK] with its corresponding prediction and the process is completed. Notably, BERT-like models such as LEGAL-BERT can yield impressive results even without fine-tuning, which means that the prediction might be accurate. This implies that occasionally, the generated chunks will match the original ones. Even though it's unlikely that this will occur for every chunk, to avoid extreme cases like this one, we handle things a bit differently.

Instead of replacing the [MASK] token with LEGAL-BERT's prediction we fetch the top 2 predictions. If the second prediction has a probability above 30%, we use the second prediction instead. If it doesn't, then we get the first one despite it having more chances of being the correct one (and hence making no changes in that chunk). An illustrative example can be seen on Figure 2.2. The percentage that we use has been fine-tuned up to a point, but in general, the logic behind this choice is that we don't want to select the second-best option if it has a great chance of being entirely wrong in terms of context and flow. Setting aside the fine-tuning, looking at the results of the predictions, the model makes correct predictions quite frequently, so we consider the approach we follow reasonable. We tested it in several instances and concluded that the changes in the resulting text were satisfactory. Finally, after every mask token is replaced by the model's predictions, the sentences are put back together and the final assembled text is an augmented version of the original one.

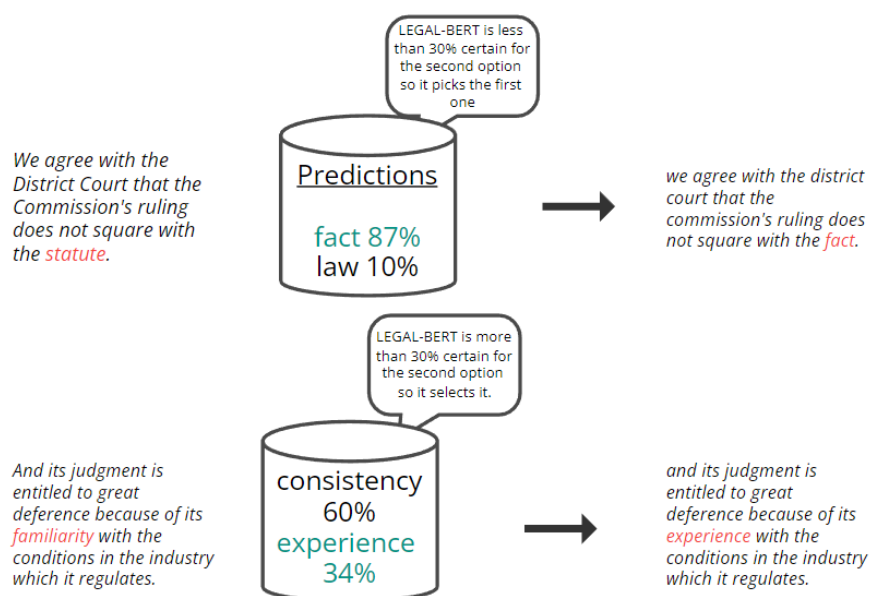


Fig. 2.2: Two different cases of masked-lm predictions applied in a chunk of our SCOTUS legal texts.

2.2.1.3 Back-Translation

The third augmentation technique is back-translation. The idea in this case is that Google Translate will translate the text from the source language which is English, to another language, and subsequently translate back to English. Google Translate is known to provide translations that preserve the overall meaning of the input most of the time, but tend to slightly change some parts. Hence, we employ Google's machine to do back-to-back translations and we expect in return altered chunks that have preserved the overall meaning.

Similar to the previous technique we described, masked-lm, we use a shifting window and break the text into chunks. This is necessary because Google Translator does not allow the input texts to be more than 5,000 characters. We will not explain the process as it is identical. The only difference is the limitation, as LEGAL-BERT required at most 512 tokens while now the issue involves the number of characters. Therefore, instead of using a shifting window of BPEs, we will implement a shifting window based on character count.

An interesting note is that we set a limit of 4,000 characters instead of 5,000. Translating from the source language to the pivot language may significantly increase the character count, potentially leading to errors in the subsequent translation from the pivot language to the destination language. This occurs because some languages demand more words than others to articulate the same information. Hence, we err on the side of caution and set the limit at 4,000.

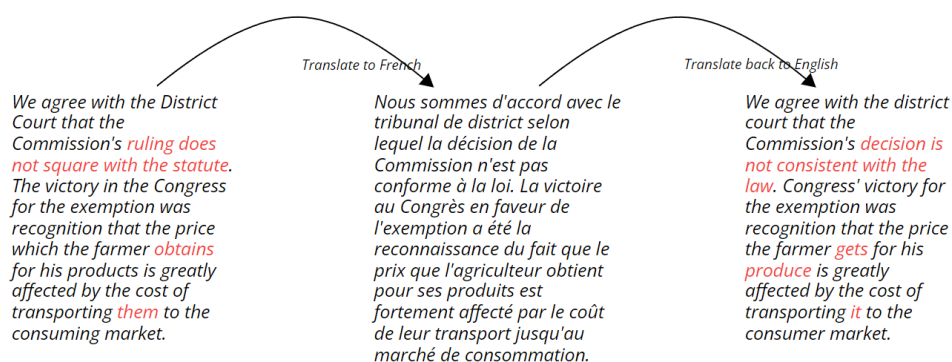


Fig. 2.3: Application of Back-Translation employing French, using a chunk of our SCOTUS legal texts. It is apparent that after back-translation has finished, a new artificial chunk has been generated. Also, the simplicity of the new chunk's vocabulary in comparison to the original chunk is noteworthy.

2.2.2 Baselines

Having introduced the three augmentation methods that will be used throughout this thesis, we will now present the baselines that we will compare them to.

2.2.2.1 Naive Oversampling Baseline

In this baseline, under-represented classes are either augmented twice, or up to the point where they have the same number of instances as the class with the most instances (N instances). This means that the majority class remains unaltered. In other words, while processing a data row, we check its label and determine if it requires augmentation. If the class has fewer than N instances, we augment it.

For example, in SCOTUS we have the following 13 classes: {'0': 1011, '1': 811, '2': 423, '3': 193, '4': 45, '5': 35, '6': 255, '7': 1043, '8': 717, '9': 191, '10': 53, '11': 220, '12': 3} and in total 5000 instances. Applying this augmentation method to SCOTUS would result in a total of 8,200 instances. This is because we set $N=1043$ and enforce each data row to be augmented up to two times, as long as the threshold N is not surpassed. The classes '0', '1', '5', and '2' will at some point reach the threshold and halt the production of augmented data rows, class '7' as we mentioned will not be augmented, and the rest of the classes will triple in size as they will not reach N . After this process is finished SCOTUS has 8,200 instances where 3,200 are artificially generated.

2.2.2.2 Masked-LM without Substitution

The masked-lm without substitution baseline does exactly what its name infers. Using LEGAL-BERT in the same manner we did for the masked-lm method (where we substitute the given token) described in [Section 2.2.1.2](#), we replace a random token with [MASK], but we never generate a predicted token that will replace the original one. So, we create one artificial instance for each data row, duplicating this way the total data rows. Ultimately, the only difference between augmented and original instances, is that the augmented ones have one [MASK] per chunk in random positions in the text, replacing the corresponding words.

2.2.3 Data Augmentation in Legal Texts

As mentioned in the previous section, prior research has laid the groundwork for data augmentation in legal texts. Since this section is aimed specifically at data augmentation methods, we will enumerate and briefly describe various augmentation techniques from recent years.

2.2.3.1 Random Deletion/Insertion/Scrambling/Swapping

These data augmentation techniques represent a fundamental approach to enhancing textual data and they are very similar in a way, so we present them unified. Random deletion involves removing random words from the text, which causes variations compared to the original text. This proves that random deletion is a valid augmentation technique, as the augmented text must be different than the original one. Wei et al. [WZ19] applied this method, resulting in slight enhancements, particularly when employing a low word removal parameter. This parameter was responsible for deciding the number of words that would be deleted from each chunk of the text. The authors concluded that too many deletions can eventually lead to text lacking coherence or meaning.

Similarly, random deletion was also applied by Santosh et al. [SBG23] as they claim that it helps the model gain a deeper contextual comprehension of the sentence, rather than depending solely on superficial word-level characteristics. The paper tackles the Rhetorical Role Labeling of legal documents, which concerns dividing a document into semantically coherent chunks and assigning a label to each chunk that depicts its corresponding role in the legal discourse. Moreover, the authors of Santosh et al. seem to agree that increasing the amount of deletions would lead to unstable results and therefore set the max deletions to be 20%. The outcome of their experiments is positive but marginal.

Furthermore, another technique is random insertion. In this approach, the authors of Wei et al. [WZ19] search for synonyms for every word in the text. A random word with at least one synonym is then chosen and one of its synonyms (if there are many) is inserted randomly into the sentence. It's essential to note that this approach may face greater challenges with more complex legal data due to their inherent complexity, which can make it difficult to find suitable synonyms that fit the context. Random insertion was also implemented by Yan et al. [Yan+19], where the task is to perform classification for crime prediction based on case descriptions. The authors of Yan et al. [Yan+19] follow a similar philosophy to Wei et al. [WZ19] but instead of words they insert whole sentences. To do that effectively, they make sure to divide the training data into clusters according to their labels, as they notice a similarity between the texts that share the same label. Afterward, they randomly pick sentences and insert them in data rows of the same label.

Additionally, the authors of Yan et al. put into effect two more augmentation methods. The second method they introduce is random deletion of entire sentences. They claim that the data have many statements that are not mandatory in their texts, and choose to delete one sentence randomly. They repeat this process for every sample and generate new instances previously unseen in the original dataset.

The third method they use is random scrambling. The authors observe that the order the sentences are fed to the model does not matter in their case. So, to introduce variety

and randomness into the dataset without changing the essential content, they perform a random scrambling operation. This shuffles the sentences within each case description text randomly, while maintaining the labels unchanged.

Finally, random swapping is a technique implemented by Santosh et al. [SBG23] which aims to swap sentences contained in a text but does that in a manner that maintains the readability for humans. They do this by limiting the exchanges of sentences with others in the same section, maintaining this way the total flow of the document. Although this might introduce some noise, the content remains intact, and the roles of sentences don't shift. This enables the model to understand how the discussion flows throughout the document and helps it handle the challenge mentioned earlier, where changes only occur within the same sections.

2.2.3.2 Data Augmentation using Semantic Relations

In their paper, Aoki et al. [AYS22] address the challenge of legal textual entailment in the domain of Japanese bar exam questions and law articles. Firstly, it should be noted that in the Japanese bar examination, students must choose 'yes' or 'no' to indicate whether an answer is entailed from a given article. The authors emphasize the need for artificial training data in this domain and propose a data augmentation method inspired by the previous works of Min et al. [Min+20] and Evans et al. [Eva+18].

The augmentation process involves identifying sentences within legal articles that contain juridical decisions or lists of juridical conditions. These sentences are split and augmented using information from previous sentences, leveraging semantic relationships. A simple example for those unfamiliar with the legal terminology, a juridical decision could be a court ruling stating that a contract is valid, while a list of juridical conditions might include specific criteria that must be met for a contract to be considered legally binding, such as the agreement of all parties involved.

Specifically, sentences containing juridical decisions are expanded with information from the preceding sentence, while sentences referring to lists of conditions are used to augment sentences discussing those conditions. This is illustrated in Figure 2.4. The resulting augmented sentences serve two purposes in the creation of artificial training data. Firstly, they contribute to the generation of positive question-article pairs, representing instances where a given statement is validly entailed by the relevant law article. In simpler words, they augment the positive sentences where the correct answer is 'yes'. On the other hand, for negative examples, sentences that involve a reversal or contradiction of the juridical decision are selected for the question parts, thereby augmenting sentences where the answer to the entailment is 'no'.

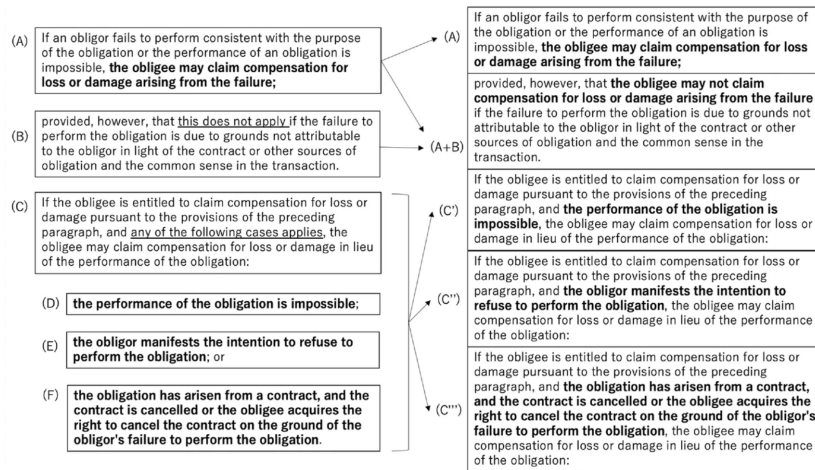


Fig. 2.4: Examples of generating sentences for data augmentation. Figure copied from Aoki et al. [AYS22]

2.2.4 Additional Data Augmentation Methods in NLP

When it comes to augmentation methods in the broader area of NLP, there are so many methods that it is not possible to mention every one of them. Feng et al. [Fen+21] kept track of a respected amount of work in this field and is a good source to get a generic perspective of data augmentation techniques in NLP. In conclusion, in the discussion on data augmentation methods, we will further mention a few additional methods. These haven't necessarily been applied in the legal or biomedical domains but are considered important for various reasons.

2.2.4.1 Synonym Replacement

Synonym replacement is an early technique that has been discussed by several research papers over the years. The idea proposed by Zhang et al. [ZZL16] was to choose random words and then retrieve their synonyms, to replace them and produce artificial data. They used an English thesaurus from WordNet [Fel05], where synonyms for a word or phrase are organized based on their semantic similarity to the most frequently encountered or widely accepted meaning. In other words, synonyms are ranked based on how closely they align with the commonly understood interpretation of the word or phrase. To determine the number of word replacements, they identified all replaceable words in the provided text and selected a random subset of size r for replacement, where r followed the geometric distribution.

This technique was some years later suggested as an augmentation method in legal document classification by Csányi et al. [CO21], but was not implemented on their behalf. It is discussed in several papers like Yan et al. [Yan+19] as well but again they chose not to

put it into effect. The consensus of the research community seems to support that even though it is an intuitive method for data augmentation, it has several formidable barriers with the most significant being the fact that it is difficult to find a manner to substitute words by their synonyms repeatedly, without altering the total meaning of the text. An example of this was presented in the [Introduction](#) with the sentence "The court decided to postpone the hearing for next week", which demonstrated that finding synonyms in a specific context is not an easy task.

It should be noted that Kobayashi et al. [Kob18] discussed the effect of synonym replacement and proposed a different idea called contextual augmentation, which is essentially identical to masked language modeling, a method that we analyzed extensively above. In this scheme, words are augmented with more varied words by using instead of synonyms, words that are predicted by a language model, given the context that surrounds the original - to be augmented - words. The only difference with our implementation is that the language model in Kobayashi et al. is a bidirectional Recurrent Neural Network (RNN) instead of a pre-trained transformer.

2.2.4.2 Data Augmentation using LLMs

In the introduction, we underlined the rising attention that LLMs have been enjoying for the past two years. It is natural, therefore, to reflect on the potential benefit of LLMs in data augmentation. This was considered by Dai et al. [Dai+23] who attempted to feed datasets to ChatGPT, the most famous LLM these days, to augment them. They do this using Amazon's dataset, PUBMED, and a Symptoms dataset from Kaggle². Their method initially included prompting ChatGPT for data augmentation. They delivered samples of all classes into ChatGPT and requested that it produce samples that maintain semantic consistency with the existing labeled examples. Subsequently, they trained a BERT-based sentence classifier on the total data, merging this way the initial and the generated from ChatGPT data, and evaluated the model. They also added numerous baselines, including word substitution, word embeddings, 2-shot ChatGPT, and many others.

To conclude, the results the authors present have both upsides and downsides. The positive part of the results is that AugGPT was the best-performing model on every dataset beating all the baselines. The downside is that they consider AugGPT to have limitations when it comes to difficult datasets like the biomedical or the legal ones in our case, due to how LLMs are trained.

In a nutshell, in this section, we attempted to present various augmentation methods put into effect in distinct domains. We first discussed a variety of augmentation techniques

²See <https://www.kaggle.com/datasets/paultimothymooney/medical-speech-transcription-and-intent> for more details

applied in a question-answering context using biomedical data. Then, we investigated methods applied specifically to the legal domain, starting with our implementations that will be evaluated in [Section 3](#). Lastly, we referred to a couple of augmentation techniques that concern the broader area of NLP.

Experiments

The majority of the time spent on this thesis was devoted to conducting experiments. In this section, we will explain in depth the whole process behind these experiments, providing all of the details in case someone wants to reproduce our results. In addition, we will include the background of the dataset and the model we employ in these experiments.

3.1 Experimental Setup

The initial objective of our experiments is clear. We aim to apply a subset of the augmentation methods proposed by Pappas et al. [PMA22] to the legal domain, utilizing the SCOTUS dataset from Chalkidis et al. [Cha+22]. The goal is to assess if some or all of our augmentation methods increase the F1-score that was produced by the model using the original not-augmented dataset.

3.1.1 Hardware Setup

Before delving into the details of our implementation, we provide the specifications of the hardware of the servers used for our experiments.

Server 1

- GPU: GeForce GTX TITAN X
- CPU: 16 CPU cores -> Intel(R) Xeon(R) CPU E5-2630 v3 @ 2.40GHz
- RAM: 32 GB
- OS: Ubuntu 22.04.1 LTS

Server 2

- GPU: 2 GPUS -> GeForce GTX 1080
- CPU: 12 CPU cores -> Intel(R) Core(TM) i7-5930K CPU @ 3.50GHz

- RAM: 64 GB
- OS: Ubuntu 22.04.1 LTS

3.1.2 SCOTUS

As we mentioned throughout this thesis so far, the evaluation of our augmentation techniques in the legal domain took with the help of LexGLUE [Cha+22], a collection of legal datasets. This collection was initially created to evaluate the performance of various models in NLU tasks. The complete collection, as detailed in the paper, comprises seven datasets. However, our primary focus was on SCOTUS, a dataset that has a significantly smaller number of training instances. This choice was driven by one principal factor. The fact that data augmentation is particularly effective when applied to datasets that are inherently limited in size. In contrast, datasets with an ample amount of data are expected to derive comparatively less benefit from augmentation.

SCOTUS stands for US Supreme Court, which is the highest federal court in the United States. It is responsible for interpreting federal law as well as the Constitution. It typically hears cases that are highly controversial or involve complex legal issues that have not been adequately resolved by lower courts. The authors designed the dataset to be a combination of Supreme Court DataBase [Spa+20] along with SCOTUS opinions generated by CourtListener a free and open-source legal research website. In a few words, CourtListener provides access to a vast collection of legal opinions, court cases, and legal documents from courts across the United States. Furthermore, SCDB provides metadata about Supreme Court cases, including details like decisions, issues, and decision directions. The SCDB covers cases from 1946 up to 2020 and it must be underlined that the data are chronologically split into training (5k, 1946–1982), development (1.4k, 1982–1991), test (1.4k, 1991–2016) sets. The task is to classify every court opinion of the SCOTUS cases in the corresponding issue area it belongs. The total issue areas are 14 but in the total data we inherit (train, development, and test) there are 13 assigned labels, as there are no instances classified as label 14.

The 14 issue areas, accompanied by their assigned integer label are the following:

0. *Criminal Procedure*: Focuses on the rules governing the process by which criminal cases are investigated, prosecuted, and decided.
1. *Civil Rights*: Involves legal issues related to the rights of individuals in society, often addressing discrimination and equal protection.

2. *First Amendment*: Encompasses issues related to freedom of speech, religion, and the press, as protected by the First Amendment to the U.S. Constitution.
3. *Due Process*: Concerns the fair treatment and procedures that individuals are entitled to in legal proceedings, ensuring fundamental fairness.
4. *Privacy*: Involves legal considerations around individuals' right to privacy and protection from unwarranted intrusions.
5. *Attorneys*: Addresses issues related to the legal profession, including standards of professional conduct and responsibilities.
6. *Unions*: Focuses on legal matters involving labor unions, collective bargaining, and the rights of workers to organize.
7. *Economic Activity*: Encompasses legal issues related to business activities, trade, and commerce.
8. *Judicial Power*: Examines the authority and jurisdiction of the judiciary, including the balance of powers among branches of government.
9. *Federalism*: Involves the distribution of powers and responsibilities between the federal government and state governments.
10. *Interstate Relations*: Concerns legal matters involving interactions and relationships between states within a federal system.
11. *Federal Taxation*: Addresses legal issues related to the power of the federal government to impose taxes.
12. *Miscellaneous*: Encompasses a broad category of legal issues that do not fit neatly into other specified areas.
13. *Private Action* (does not exist in our data): Pertains to legal issues involving private individuals taking legal action against each other.

SCOTUS includes 5,000 training instances from the period 1946-1982, 1,400 evaluation instances from the following nine years (1982-1991), and finally 1,400 test instances from 1991-2016. We present a specific trait that this dataset has, class imbalance.

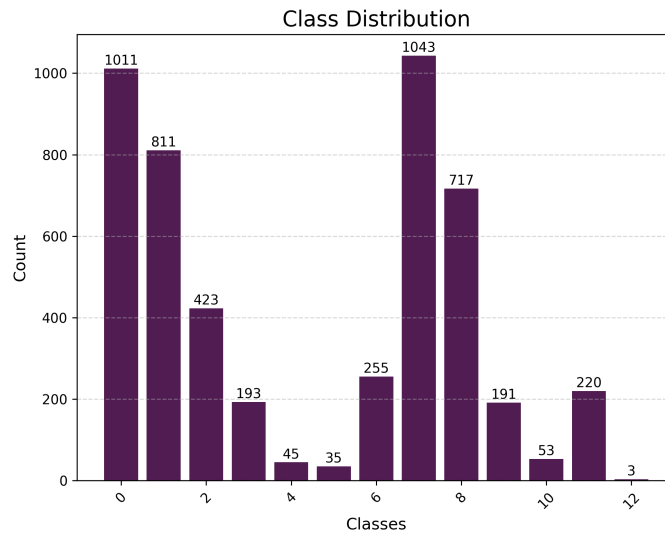


Fig. 3.1: A bar chart of the class distribution in SCOTUS

It is apparent that the classes are far from equally distributed. There is a class with only 3 instances and that is class 12 (see Section 3.2), which corresponds to the issue area called "Miscellaneous". Moreover, there is a class that corresponds to the Private Action issue area and is not represented at all in our data (it does not exist in development or test data as well).

3.1.3 LEGAL-BERT

Researchers from LexGLUE [Cha+22] and their team, employed the benchmark dataset they generated, to evaluate a variety of models. The best-performing model overall was LEGAL-BERT developed by Chalkidis et al. [Cha+20]. Hence, it is a natural choice to select this model to run our experiments.

LEGAL-BERT is a BERT-like model designed to support legal NLP research, computational law, as well as legal technologies. The authors produced a family of models, each one with different characteristics, and the model we are using is one of these variants. In a few words, the authors underline the mediocre results of BERT in scientific areas like the biomedical one or the legal domain, despite the fine-tuning of relevant datasets. Therefore they conduct several experiments and propose two main methods, further pre-training a BERT on legal corpora and pre-training a BERT on legal corpora. They employed three datasets, which we will briefly discuss below, and concluded that the best-performing model was BERT pre-trained from scratch on all three datasets combined, outperforming BERT pre-trained from scratch on each dataset separately, as well as each variant of further pre-trained BERT (regardless of the datasets that it was pre-trained).

The datasets that were employed for the pre-training of LEGAL-BERT were EURLEX57K, ECHR-CASES, and CONTRACTS-NER. EURLEX57K, a dataset introduced by Chalkidis et al. [Cha+19] is a multi-label classification dataset that encompasses European Union laws. In addition, ECHR-CASES is a multi-label classification dataset (also suitable for binary classification) that involves cases from the European Court of Human Rights. Lastly, CONTRACTS-NER contains US contracts. It is a dataset that was developed for named entity recognition and consists of three separate parts that center around contract headers, dispute resolution, and lease details.

Subsequently, we will compare LEGAL-BERT and BERT to assess the impact of the modifications introduced by the authors, that resulted in significant improvements. Both LEGAL-BERT and BERT share the same architecture, featuring 12 Transformer encoder blocks, 768 hidden units, and 12 attention heads, and have been trained with 110 million parameters. This completes the creation of a new vocabulary of the same size as BERT's vocabulary.

LEGAL-BERT underwent training for 1 million steps, roughly equivalent to 40 epochs, using all three datasets mentioned earlier. The training was conducted in batches of 256 samples, with sequences of up to 512 sentence-piece tokens. Additionally, the authors utilized the Adam optimizer with a learning rate set at $1e-4$, mirroring BERT's configuration. Hardware-wise, the authors performed training using the official BERT code on v3 TPUs equipped with 8 cores from Google Cloud Compute Services. In terms of tuning, they experimented with batch sizes from the set 4, 8, 16, 32. Furthermore, they made adjustments, such as lowering the learning rate to $1e-5$ to prevent overshooting local minima and increasing the dropout rate to 0.2 to enhance regularization.

The results achieved by fine-tuning LEGAL-BERT, which we adopted in our experiments, consistently outperformed fine-tuned BERT models. Notably, the most significant improvements were observed in multi-label tasks, where the LEGAL-BERT variations showcased remarkable enhancements.

3.1.4 Hierarchical Model

LEGAL-BERT is a model that is expected to perform particularly well in a classification task where the corresponding dataset is relevant to legal issues. And yet, LEGAL-BERT that was analyzed in the previous section, does not include any implementation that would make it competent to read and classify long texts, like the ones our dataset includes. For this reason, when the authors of LexGLUE evaluated a variety of models using these datasets, they created a script that "hacks" the BERT-like model's implementation to add a hierarchical variant introduced in Chalkidis et al. [Cha+21].

The hierarchical variant addresses the challenge of processing long texts efficiently. It leverages pre-trained Transformer-based models like BERT in our case, to independently encode each paragraph in the input text. This process results in paragraph representations (each one being the output embedding of the [CLS] pseudo-token, which represents the entire input, in our case a paragraph), that hold essential information. To make these paragraph representations context-aware, the authors employ a second-level shallow Transformer encoder. This encoder maintains consistent specifications like hidden units and number of attention heads, across different pre-trained models and ensures that the paragraph representations are aware of the surrounding paragraphs. The hierarchical variant combines these context-aware paragraph representations through max-pooling and yields a document representation, which is then passed through a classification layer for downstream tasks like text classification in our case.

The process described above is the exact implementation that takes place in the code we inherited from Chalkidis et al. [Cha+22]. We used LEGAL-BERT as our base encoder and then we called the hierarchical encoder which comprises two Transformer layers. The first one encodes paragraphs independently, and the second ensures context awareness. During the forward pass (the process where input data is fed into LEGAL-BERT), the input is split into paragraphs before it's encoded and transformed into a document representation. We then used this representation as a foundation for the classification that happens in the next step.

The authors of LexGLUE, for datasets like SCOTUS that have long enough texts to need the hierarchical variant, follow a specific approach. Here, they take the document representation obtained by max-pooling across paragraphs. This representation is then passed through a linear layer for classification. Depending on the dataset, this linear layer is followed by a softmax activation. SCOTUS is designed as a multi-class classification dataset, where each input belongs to one and only one class, making Softmax a standard choice for this scenario.

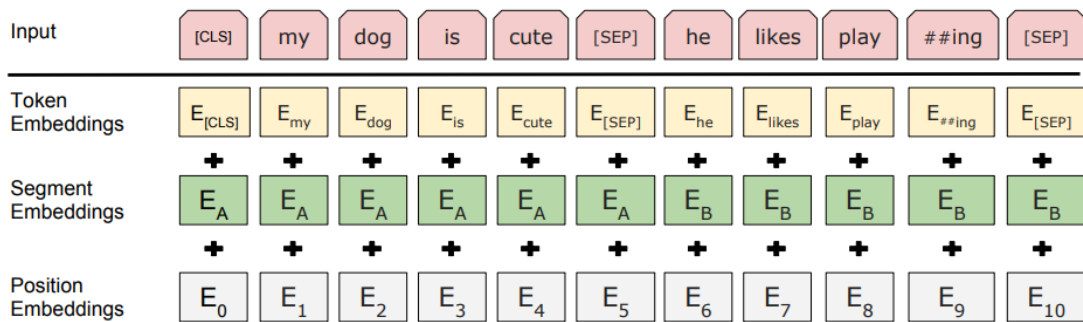


Fig. 3.2: BERT input representation. The input embeddings are the sum of the token embeddings, the segmentation embeddings, and the position embeddings. In the final hidden state (when the second [CLS] is processed), the model captures the overall information from the entire sequence. Figure copied from Devlin et al. [Dev+18]

In addition, at this point, we should state that all of the experiments that we conducted have one important thing in common. Due to limited computing power, during the training, we froze all of the layers except the last one which handles the long training instances and classifies the text. This means that before we started developing the augmentation methods, we needed to run the frozen model and feed the non-augmented data (the original data of LexGLUE) to get a fair baseline.

Also, due to the results being unstable and occasionally causing the model to halt training, we enforced gradient clipping, which is a technique that prevents exploding gradients [PMB13]. Exploding gradients occur when the gradients of the model’s parameters become extremely large, causing training instability and slow convergence. Gradient clipping, is a technique that limits the size of gradients during training to a specified threshold, preventing them from becoming too large. In our case we set the maximum norm to be equal to 1, hoping this way to restrict the unstable nature of the results. In general, this value should have been fine-tuned or at least have been selected with a sufficient explanation regarding the choice. We set the maximum norm to be equal to 1 without any fine-tuning or sufficient arguments to justify it. The reason is that we aim to repeat our experiments without layer freezing or gradient clipping, using exactly the settings of LexGLUE, when there is enough time and computing power available. As far as this thesis is concerned, we focus on the comparison between augmentation methods and a selection of baselines and ignore the bigger picture, which is to find the best methods and optimize their hyper-parameters, achieving this way the best possible results. Moreover, it should be mentioned that when we run the frozen model and feed the non-augmented data of LexGLUE, we also enforce gradient clipping for the same reason as before, to make a fair comparison. The results of this run will constitute the baseline for our experiments.

3.1.5 Evaluation Metrics

The F1-score [SJS06] is a popular evaluation measure in classification tasks and is particularly useful when dealing with imbalanced datasets. Accuracy, which is the most common metric, cannot be trusted for imbalanced datasets, since a baseline model that classifies all instances to the over-represented class, achieves high accuracy scores. F1-score avoids that, as it has two other terms that handle things differently as its foundation, precision, and recall. Before continuing our analysis regarding the F1-score, let’s quickly define its building blocks: Precision and Recall.

$$Precision = \frac{TP}{TP + FP}$$

TP, or True Positive, denotes instances that are correctly predicted to belong to the positive class. In simpler terms, it represents the number of true positive examples accurately identified by the model. Similarly, FP, or False Positive, refers to instances that are incorrectly predicted as belonging to the positive class. This category contains examples that the model mistakenly identifies as positive.

$$Recall = \frac{TP}{TP + FN}$$

FN, or False Negative, indicates instances that are misclassified as negative by the model. In other words, these are cases where the model fails to identify positive examples. It's important to note that True Positive, False Positive, and False Negative are metrics calculated per class or category.

It should be noted that precision and recall are performance metrics computed per class. Precision measures the accuracy of the model when it predicts a positive class, while recall assesses the model's ability to capture all positive instances.

Having covered the foundation of the F1-score, we can now define the F1-score to be the harmonic mean of precision and recall. It balances the trade-off between precision and recall and provides a single score that reflects a model's overall performance. Its formula is the following:

$$F1 = \frac{2 \cdot precision \cdot recall}{precision + recall}$$

Moreover, in an attempt to understand intuitively the F1-score's functionality, a high F1-score means that precision and recall both have high values. This shows that the model is making accurate positive predictions and it shifts away from generating many false negatives. On the other hand, a low F1-score is an indication of an imbalance between precision and recall, meaning that conversely to the previous case, the model is either making too many false positive errors or missing many positive instances. This shows that the F1-score is a quality metric when we seek a balance of precision in recall, and we almost always do. By letting it steer our evaluation measures, we gain a solid understanding of our model's performance and that is a crucial reason why some people prefer it over accuracy, which has the advantage that it is very intuitive, but it is also prone to be misleading.

Having settled on the aspects of F1-score it's high time we take it a step further. The main metric that we consult to monitor the progress of our model is indeed the F1-score but with a slightly altered scheme. In our case in SCOTUS, we have multi-class classification,

not binary classification. In this case, the F1-score has to adapt its calculations due to the greatest amount of classes and there are two different ways to handle that.

The first case is to consider each class independently and then take the average of these F1-scores. This is done by initially calculating the F1-score for each class separately using the standard formula for binary classification and then calculating the average F1-score. In this case, F1-score is called "Macro F1-score" and is defined by the formula:

Macro F1-Score:

$$\text{Macro-F1} = \frac{1}{N} \sum_{i=1}^N \text{F1}_i$$

The second case is to consider all instances across all classes as a single large binary classification problem. This happens by first computing the total number of True Positives (TP), False Positives (FP), and False Negatives (FN) over all classes and subsequently calculating micro-precision and micro-recall using the aggregated values. In this case, the F1-score is called the "Micro F1-score" and its formula is:

Micro F1-Score:

$$\text{Micro-F1} = \frac{2 \cdot \text{Micro-Precision} \cdot \text{Micro-Recall}}{\text{Micro-Precision} + \text{Micro-Recall}}$$

3.2 Experimental Results

We assess our models in the same way that Chalkidis et al. [Cha+22] did, to compare fairly the results. For each dataset and each augmentation method, we run LEGAL-BERT with 5 distinct seeds. The final results are generated by calculating on the test set the mean of the 3 best-performing seeds selected based on the macro F1-score of the development dataset, which is exactly how Chalkidis et al. produce their results.

3.2.1 SCOTUS Results

3.2.1.1 Basic & Smart Augmentation

Moving on, we present the results of the inherited data where augmentation has no involvement.

The last column depicts the final results of LEGAL-BERT using the original data. It shows the average of seeds 2, 3, and 5 that had the best macro-F1 performance on the development dataset. These results will be our main baseline for the rest of the experiments.

	Seed 1	Seed 2	Seed 3	Seed 4	Seed 5	Results
dev loss	0.796	0.856	0.899	0.778	0.767	0.841
dev macro-F1	0.715	0.737	0.740	0.722	0.728	0.735
dev micro-F1	0.797	0.809	0.812	0.805	0.806	0.809
test loss	0.952	1.081	1.257	0.976	1.031	1.123
test macro-F1	0.665	0.694	0.645	0.684	0.669	0.669
test micro-F1	0.759	0.773	0.749	0.764	0.762	0.761

Tab. 3.1: LEGAL-BERT's performance using the original data for 5 different seeds. The last column "Results" takes the average of the 3 best-performing seeds. The best-performing seeds are determined by their results in the category of development macro-F1. The bold cells indicate the best-performing seeds and the total results.

We continue with the next two baselines that we defined a few lines above. We will not present again the results of each seed from this moment forward, every result that we introduce has been generated with the same methodology. The runs of the two baselines can be seen at [Table 3.2](#). We notice that the two baselines that we defined are weaker than the data without augmentation. For the latter method, we ran the model with 5,000 training instances, whereas for the two baseline methods, we had more due to augmentation. To figure out the number of training instances we need to go back a few pages where we described how they work. Starting from masked-lm without substitution see [Section 2.2.2.2](#), the total number of instances is 10,000, as for every data row we created a new one. Furthermore, naive oversampling contains 8,200 examples and we can compute that following the calculations we provided in [Section 2.2.2.1](#) when defining this baseline.

At [Table 3.2](#) we also introduce the first augmentation technique, word2vec. The first experiment we conducted had a similarity equal to 0.95 and the available augmentations per row were set to 1. This essentially means that this run fed LEGAL-BERT with the 5,000 original data rows plus at most 5,000 artificial data rows. We underline "at most" as it is possible that some data rows were not available for augmentation, due to the embedding search failing to find neighbors that fulfill the similarity threshold. In this specific case though, we report that we generated exactly 5,000 artificial data rows.

	Original Data	Naive Oversampling	Masked-LM w/out Subst.	Word2vec 0.95
dev loss	0.841	0.879	0.856	0.770
dev macro-F1	0.735	0.725	0.73	0.733
dev micro-F1	0.809	0.809	0.803	0.808
test loss	1.123	1.281	1.152	1.010
test macro-F1	0.669	0.647	0.659	0.674
test micro-F1	0.761	0.748	0.755	0.763

Tab. 3.2: Data Augmentation using Word2vec compared to the 3 baselines. The bold cells indicate the best-performing method in each category. Word2vec with the current settings (0.95 similarity threshold), appears to surpass the baselines in the test set but fails to do so in the development set.

After this run, we question whether increasing the number of augmented instances (for example we could generate up to 2 from each data row), or lowering the similarity threshold would produce better results. Due to time and computing power restrictions, instead of conducting similar experiments with 5 seeds for each hyperparameter tweak, we selected the best-performing seed for word2vec (which was seed 3) and tried different settings running only this seed. The results of this process are portrayed below. The first column represents the augmentation we performed so far, augmenting each data row once with the similarity threshold set at 0.95. The second column is identical to the first one with one important difference, instead of augmenting one time each data row, we augment it twice. Finally, the last two columns augment each data row once but test various similarity thresholds.

	Word2vec 0.95	Word2vec 0.95, x3	Word2vec 0.9	Word2vec 0.85
dev loss	0.726	0.896	0.791	0.945
dev macro-F1	0.723	0.709	0.733	0.720
dev micro-F1	0.808	0.795	0.81	0.803
test loss	0.880	1.184	1.070	1.296
test macro-F1	0.670	0.668	0.672	0.615
test micro-F1	0.762	0.750	0.764	0.716

Tab. 3.3: Word2vec hyper-parameter tuning. The bold cells indicate the best-performing method in each category. All of these results were produced by seed #3 which was the best-performing seed. In this table, we compare different similarity thresholds as well as different amounts of generated artificial instances (as the second column involves tripling the data). Word2vec with 0.9 and 0.95 similarity thresholds appear to be the best options. "Dev" loss, corresponds to the loss of the development dataset.

These results show that tripling the data is not beneficial, so from this point forward, we chose to augment each data row at most 1 time when applying word2vec. Additionally, even though the performance of word2vec using 0.85 as a threshold was much worse than the baselines, word2vec with 0.9 as a threshold was inspiring enough to motivate us to run all of the seeds. We present the comparison of word2vec with a 0.95 threshold versus 0.9 below.

	Word2vec 0.95	Word2vec 0.9
dev loss	0.770	0.816
dev macro-F1	0.733	0.740
dev micro-F1	0.808	0.808
test loss	1.010	1.054
test macro-F1	0.674	0.666
test micro-F1	0.763	0.760

Tab. 3.4: A head-to-head comparison of Word2vec with different similarities. The generated results depict the average performance of all 5 seeds (following the same notion as in all the previously reported results). Word2vec with 0.9 similarity is the best method for the development data. The bold cells indicate the best-performing method in each category.

The head-to-head comparison shows that word2vec with 0.9 similarity is slightly better on the development data, but worse on the test data. To avoid overfitting the development data we chose word2vec with 0.9 similarity as our main word2vec implementation.

We proceed to the subsequent augmentation method, back-translation. We conformed with the implementation of Pappas et al. [PMA22], who initially tried back-translation with only French as the pivot language. Afterwards, we conducted the same experiment using French, German, and Spanish as pivot languages which means that this implementation has 20,000 instances in total. We present the head-to-head comparison between these back-translation implementations.

	Back-Translation /w French	Back-Translation /w all Languages
dev loss	0.959	1.361
dev macro-F1	0.734	0.730
dev micro-F1	0.806	0.806
test loss	1.294	1.830
test macro-F1	0.668	0.656
test micro-F1	0.761	0.749

Tab. 3.5: A head-to-head comparison of Back-Translation with different languages. The bold cells indicate the best-performing method in each category. Back-translation using only French as a pivot language beats back-translation that employed 3 languages (French included) and becomes the default back-translation for the remainder of our experiments.

We conclude that back-translation with only 1 pivot language is the best option as it beats back-translation with 4 languages in almost every metric. Hence, after reviewing these results we held this as our back-translation method.

Finally, we put into effect our final augmentation method, masked-LM. We evaluated our model with 10,000 instances, where in each augmented row one word has been predicted by LEGAL-BERT. The table below compares masked-lm with the previous two augmentation methods, along with the three baselines we defined. It is a total comparison between all of the augmentation methods we applied and all of the baselines.

	Original Data	Naive Oversampling	Masked-LM w/out Sub.	Word2vec	Back-Translation	Masked-LM
dev loss	0.841	0.856	0.879	0.816	0.959	0.780
dev macro-F1	0.735	0.730	0.725	0.740	0.734	0.732
dev micro-F1	0.809	0.803	0.809	0.808	0.806	0.808
test loss	1.123	1.152	1.281	1.054	1.294	1.018
test macro-F1	0.669	0.659	0.647	0.666	0.668	0.666
test micro-F1	0.761	0.755	0.748	0.760	0.761	0.764

Tab. 3.6: A full comparison between all methods and baselines across all seeds. The bold cells indicate the best-performing method in each category.

Data augmentation was successful regardless of the method used in Pappas et al. [PMA22] where we inherited the methods from. The differences between the results of the methods concerned the scale of the benefit each method offered. The results generated on the legal domain, for the SCOTUS dataset, which are depicted in the above table, do not point in the same direction. None of the augmentation methods managed to surpass the performance of the initial data. In the subsequent part of this section, we seek potential reasons that might have caused this.

The unexpected results show that augmentation not only does not benefit, but it leads to worse results, which prompted a different idea. If the class imbalance is indeed affecting negatively the model's performance, then by augmenting all of the data rows the results may be also affected negatively. For instance, the majority label that has 1,043 instances will have 2,086 in the augmented dataset, but the minority class with initially 3 instances will only have 6.

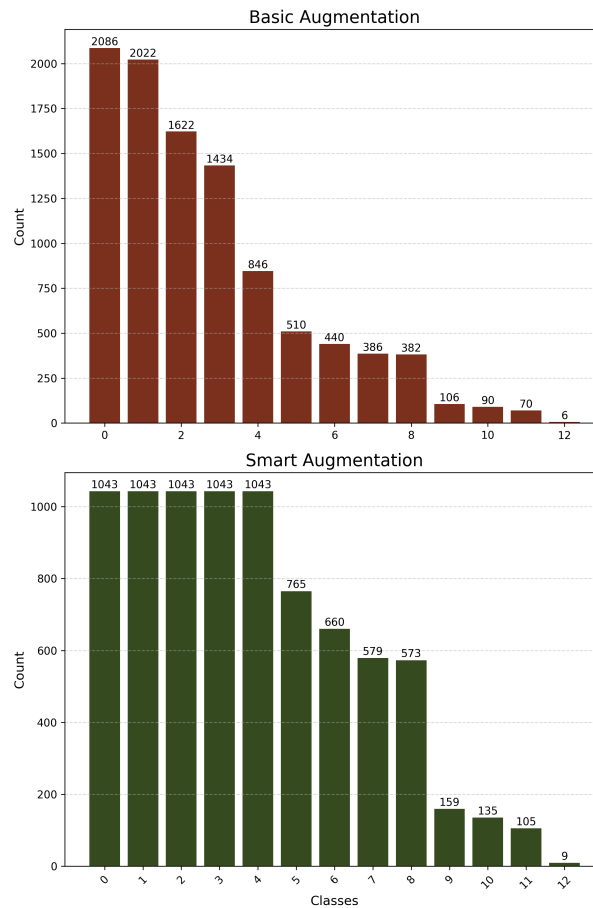


Fig. 3.3: A bar chart class distributions in SCOTUS. The first figure shows the class distribution using the basic augmentation method, whilst the latter demonstrates the class distribution using "smart" augmentation.

To address the problem of class imbalance, we conduct all of our experiments again (for the 3 augmentation methods) using a different technique than the one described in the previous section. Instead of doing what we did so far, which was essentially generating one artificial data row without accounting for its label, we augment rows, in the same manner we augmented duplicate data in naive oversampling. So, we only augment the rows that are under-represented, and make sure not to surpass the threshold N , which corresponds to the number of instances of the majority class. Therefore, our dataset will now contain at most 8,200 instances, which is the sum of the numbers shown in each column of the second figure in Figure 3.3. In this figure, we can see the difference in the class distribution between the two augmentation techniques. From this point forward we will call the augmentation method that we applied so far as basic or standard augmentation, and the alternative technique we propose will be called "smart" augmentation.

Ultimately, we present the results of this alternative augmentation approach. The generated results have been produced with the same circumstances used previously for the standard augmentation method, to make a valid comparison.

	Original Data	Naive Oversampling	Masked-LM w/out Sub.	Word2vec	Back-Translation	Masked-LM
dev loss	0.841	0.856	0.879	0.935	1.008	0.859
dev macro-F1	0.735	0.730	0.725	0.741	0.733	0.730
dev micro-F1	0.809	0.803	0.809	0.807	0.800	0.805
test loss	1.123	1.152	1.281	1.284	1.329	1.130
test macro-F1	0.669	0.659	0.647	0.659	0.660	0.654
test micro-F1	0.761	0.755	0.748	0.750	0.748	0.752

Tab. 3.7: A full comparison between all methods and baselines, across all seeds, using smart Augmentation. The bold cells indicate the best-performing method in each category. We observe that the results of the original data remain the best overall, whilst the margin between them and the 3 augmentation methods increases.

Analyzing the data presented in the table above, it becomes evident that none of the employed augmentation methods managed to improve the results achieved by Chalkidis et al. [Cha+22]. Word2vec surpassed the original data in development macro-F1, but regarding the test results, each augmentation method failed to score a better result. Surprisingly, the more sophisticated augmentation approach that we expected to improve the results, performed even worse than the standard one. To shed light on these perplexing scores, we will employ data visualization in the form of figures. These visual aids will enable us to assess the model's learning progress across different augmentation techniques and will hopefully explain the background of the generated results. Furthermore, we will shift our focus to the confusion matrices, along with the F1-scores for each class, to get a better understanding of the performance of each distinct augmentation method.

3.2.1.2 Visualization & Analysis

To begin, we will examine the training and development losses associated with each augmentation method. Our approach involves selecting the best and worst-performing seeds for each method and plotting their results on a shared figure or sub-figure, as appropriate. In total, we will generate four plots, one for the original non-augmented data, one for back-translation, one for masked-lm, and one for word2vec. Each of these plots will contain four distinct curves: the training and development curves for both the best seed and the worst seed.

This method of analysis allows us to compare the performance of the best and worst seeds and gain valuable insights into the learning process for all the augmentation techniques. It's important to note that the curves for different seeds will conclude at varying epochs, as they correspond to a different run.

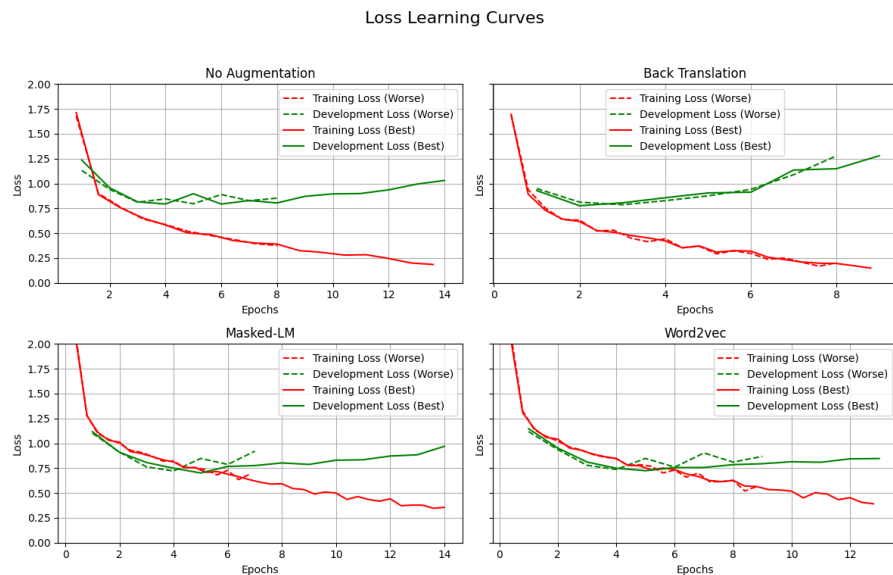


Fig. 3.4: Loss Learning Curves for the best and worse seed of each method. We can see that seeds within each method display similar trends, with word2vec having the largest margins. Overall, the model's performance seems to achieve better results in masked-LM and word2vec, achieving less overfitting. The loss learning curves are derived from the optimal model versions, specifically the ones from the 4th epoch counted in reverse. This selection is guided by the application of early stopping after 3 epochs, ensuring that subsequent epochs do not contribute to improvement.

The first impression of the learning curves is that the model is performing well. We achieve a lower validation loss with word2vec and masked-LM, while the gap between training and development loss is smaller, and the development loss curve increases less rapidly when the overfitting starts. This is expected as we added more data and therefore achieved better regularization. Back-translation does not seem to reduce the loss, but recall that it

produced the best results in the test data, showing that once again that a smaller loss does not promise a better F1-score.

Furthermore, it is important to note that by augmenting the dataset we doubled the training data, and one epoch has twice as many update steps. This means that if for example, we do 12 epochs with word2vec, in the augmented version it's like word2vec is doing 24 epochs in terms of update steps. This means that if we hadn't performed the augmentation the development loss would have soared much higher than 1.0, but instead it declines significantly. This result is promising, and possibly additional tuning to the learning rate could lead to better results.

Regarding the differences between the seeds, word2vec has the biggest difference between the seeds while back-translation's seeds are the most alike. To provide a clearer view, the difference between the best seed and the worst seed of back-translation was roughly 1.1% whilst for word2vec it was 2.5%. These results overall show that something else is off, as the loss learning curves are generally performing well, so we investigate further the behavior of the model by generating the confusion matrices.

We begin by calculating the confusion matrix of the original data which was the run that performed best overall. Then, we do the same for the best-performing standard augmentation method and the best smart augmentation method. Back-translation had the highest macro F1-score in both cases which is our main point of comparison, so we selected it. Lastly, we need to select a seed as we have 5 seeds per method, so we choose the best-performing seed for each method.

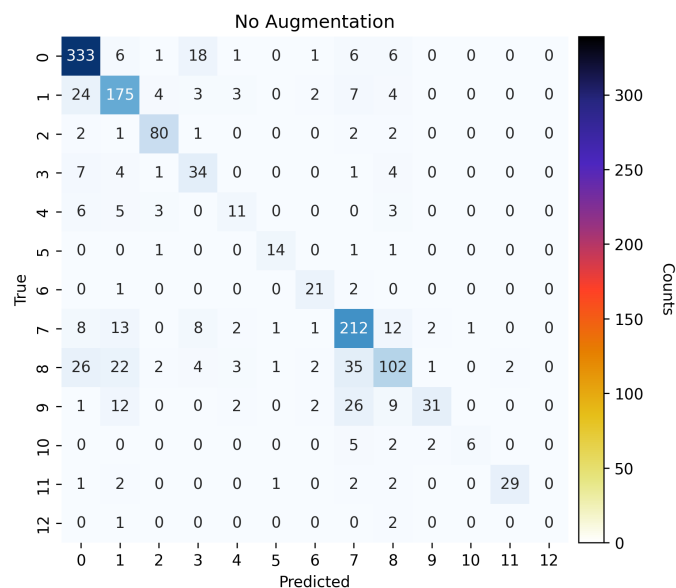


Fig. 3.5: Confusion Matrix # No Augmentation

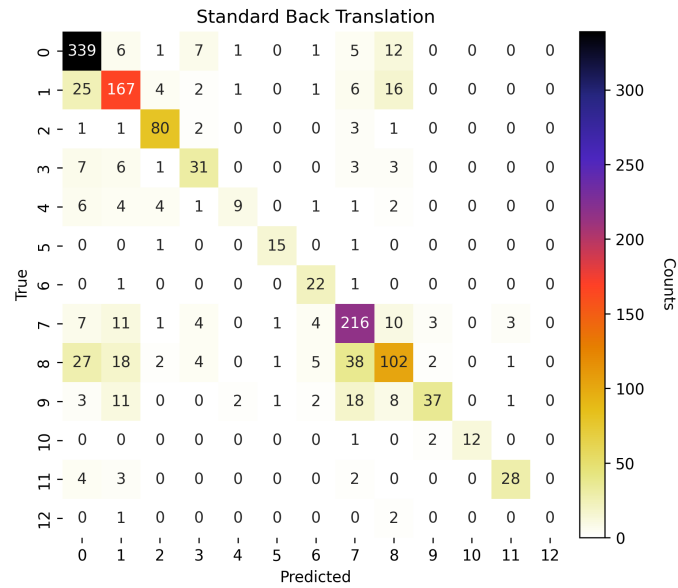


Fig. 3.6: Confusion Matrix # Standard Augmentation

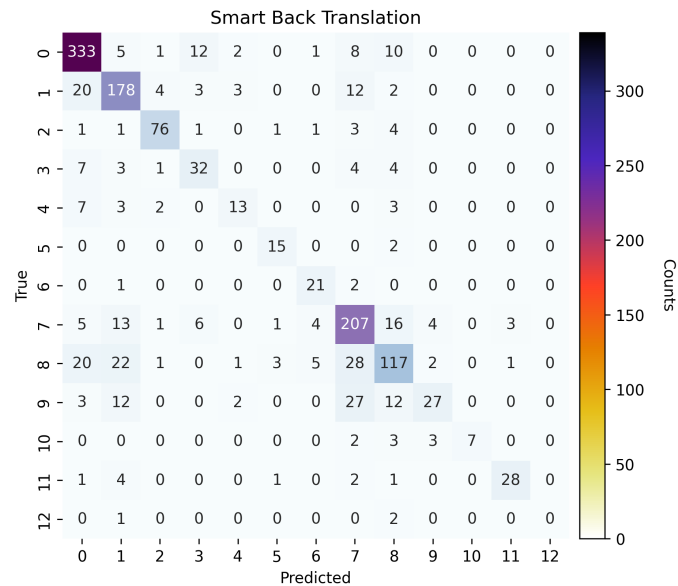


Fig. 3.7: Confusion Matrix # Smart Augmentation

What we can derive from the matrices is that class '12' has 3 instances on the test set and the model fails in every method to get one correct. This is not a surprising result as the training instances are far too low, so low that augmentation can't change the result.

Furthermore, the classes '4', and '5' which also had very few examples in comparison to others, and we expected them to demonstrate significant improvements, seem to yield improvements. Class '4' achieved 11/28 without augmentation, 9/28 with standard augmentation, and 11/28 with smart augmentation. Moreover, class '5' achieved 14/17 without

augmentation and 15/17 in both of the augmentation methods. In contrast, classes '3' and '11' which were also severely under-represented, do not benefit from any of the augmentation methods. Class '3' achieved 34/51 without augmentation, 31/51 with standard augmentation, and 32/51 with smart augmentation. Finally, class '11' achieved 29/37 without augmentation and 28/37 in both of the augmentation methods.

The above results demonstrate that it is unclear if augmentation is helpful and which of the two is best, as each label is affected differently. In addition, it should be mentioned that most of these changes are relatively small and there is a possibility that they are random.

Subsequently, we calculate the F1-score per class.

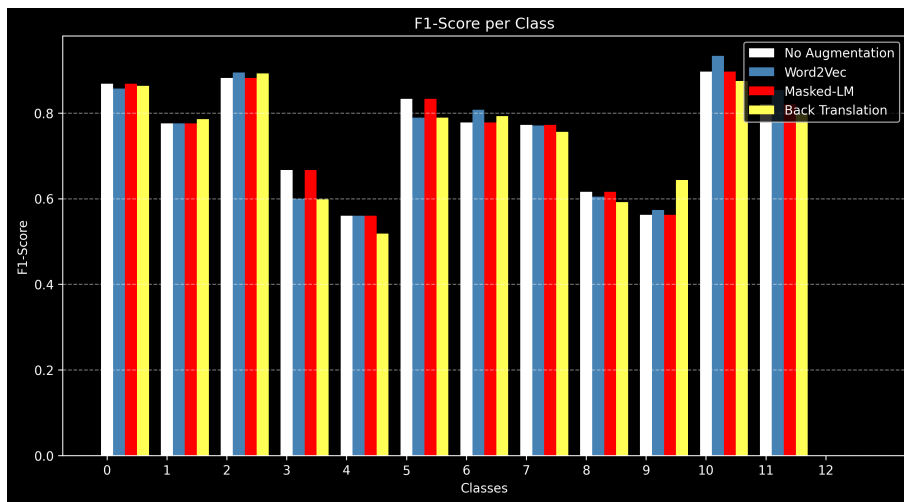


Fig. 3.8: F1 per class with standard augmentation, using the development results of our runs.

By examining the figure we see that there are shifts for each label. These shifts can be attributed to the differences we discussed previously in the confusion matrices. The results between each method are very close, so we can't be confident that there is a reason behind them. Assuming that this is not occurring randomly, it is an interesting result as it shows that each augmentation method captures different kinds of details. For example, back-translation fails to increase the performance in label '3' whatsoever, but outputs great results in comparison to the other methods for label '9'.

Then, we repeat the calculations for smart augmentation:

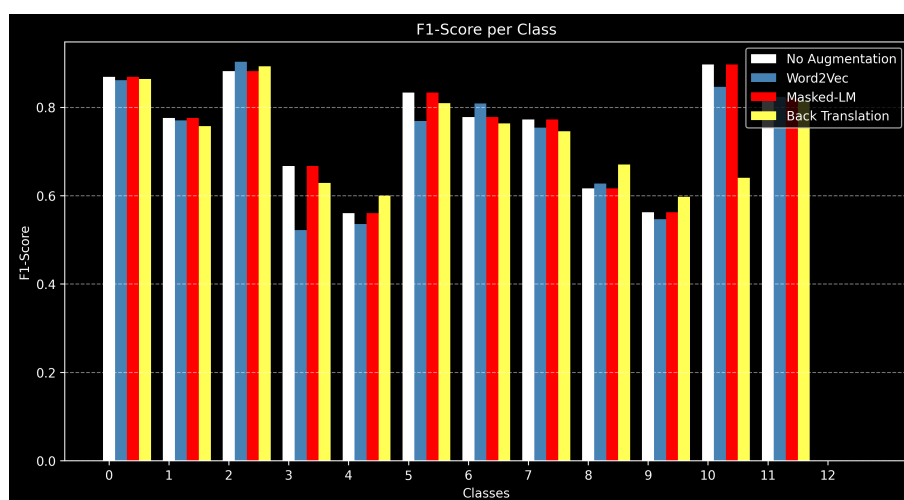


Fig. 3.9: F1 per class with smart augmentation, using the development results of our runs.

The results reported in [Figure 3.9](#) demonstrate that some labels are heavily affected by the choice of smart augmentation. The most notable shift involves label '10', as the back-translation results are much worse than those in the previous figure, showing a radical decline. This is not the case for every label though, as labels '4' and '8' using once more, back-translation show an improvement in comparison with standard augmentation.

Furthermore, we can observe additional interesting shifts in label '10', as using word2vec the performance of smart augmentation once again decreases significantly compared to standard augmentation. The perplexing observation stems from masked-lm though, as, in contrast with back-translation and word2vec, it seems to improve its results in label '10' using smart augmentation. This last note supports our previous claim that each method and each label capture different traits of the data, a very interesting result.

Overall, the differences between the F1 per class figures demonstrate that it is not clear which of the two augmentation schemes, smart and standard, is the best. Each class is affected uniquely, showing that currently, it is hard to reach a safe conclusion.

The final thoughts behind LEGAL-BERT's performance in SCOTUS with the addition of artificial data are mixed. None of the techniques achieved a rise in macro-F1 for the test data, showing that changing the dataset only leads to worse results. The tables and figures we included were also pointing in the same direction.

An interesting observation is that the more we interfere with the initial data using word2vec, the worse the results are. We might have used word2vec with a 90% similarity threshold after all, but as we presented in our results, on the test set the best method was the one with 95% as a lower bound, and even that did not manage to improve the original test macro F1-score. It is natural that the lower the threshold is, the bigger the amount of alterations

in the initial text, so this could potentially be a reason. For instance, the experiment with similarity set to 85%, which is far from low, resulted in a 5% macro F1-score dive. It is unlikely that the words replaced are so off context to explain this dive, as 85% is a relatively high threshold.

Similarly, back-translation as we explain in [Figure 2.3](#), tends to simplify complicated words during the pivot translation process, and specialized words are frequent in legal texts. This could be a potential reason why this method fails to deliver the expected results. It should be mentioned though that this method changes significantly the initial text (definitely a lot more than the other two methods), and yet produced the best macro-F1 results even if it was a marginal result. This clashes with the assumption we made regarding the negative relation between good performance and the number of word changes in the initial text.

Moreover, masked-lm is the method that is expected in theory to make the smallest amount of word changes in comparison with the other two methods. This is because it changes only 1 word per chunk, while the other two methods don't have a limit. The fact that masked-lm managed to surpass the original data baseline in three categories (development loss, test loss, and micro F1-score) supports our claim that perhaps changes are affecting negatively the results.

Conclusions and Future Work

4.1 Conclusions

In this thesis, we augmented the training instances of SCOTUS, investigating the impact of augmentation in a legal context. Having performed and visualized a number of experiments, we conclude that augmentation did not manage to succeed in our case. None of the word2vec, back translation, and masked-lm that reported optimistic results in previous work, achieved to surpass the benchmark. Even when we tried to adapt to SCOTUS's peculiarities, which had a substantial class imbalance, the outcome was the same.

The various methods we tried did not result in an improvement in the macro-F1 score for the test data. The first results we presented in [Table 3.2](#) showing the performance of the two "strong" baselines, were an early indication of that, as none of them managed to surpass the original results. Word2vec was the only method that managed to beat the best baseline in one category involving F1-score and that was the macro-F1 of the development dataset. This was also confirmed by the loss learning curves, as its development curve was notably improved in comparison with the one without augmentation. At the same time though, it performs quite worse on the test data.

Back-translation was the best technique in comparison with the other augmentation techniques, as it obtained 66% macro F1-score beating word2vec that output almost identical results (65.9%), despite depicting a higher loss. Also, masked-lm had much better validation and prediction loss than every other method which is also supported by the learning curves, but had poor performance in terms of F1-scores, as it was 1.5% below the benchmark in macro-F1 for the test data. Finally, the additional experiments we conducted applying a different style of augmentation by prioritizing the minority labels, despite the optimistic expectations proved to be even worse than our standard augmentation plan.

Altering the dataset ultimately led to poorer performance in every method, and the results of word2vec might be the most clear indication of that. By lowering the similarity threshold from 95% to 85% we observe a decline of the performance, which is noteworthy. Meanwhile, another interesting remark, is that some methods seem to perform better in separate labels, potentially demonstrating that some techniques might fit better to some classes. It has to be mentioned though that these discrepancies are quite low most of the times, and they do not constitute a concrete argument. The difference using back-translation and comparing

the results of smart and standard augmentation in label '10' though, are promising, and show that there is room for more research in this area.

4.2 Future Work

The direction that Dai et al. [Dai+23] selected, who fed datasets to ChatGPT in order to augment them is interesting (and obvious considering the hype behind LLMs nowadays), and it just might be the method that will unlock the potential of data augmentation in legal texts. Therefore, using LLMs to generate artificial data is a very intriguing idea for the legal domain, and for data augmentation methods in NLP in general. Moreover, reflecting on the results our experiments generated where there seemed to be a universal problem regardless of the augmentation technique, an interesting approach would be one suggested by Csányi et al. [CO21]. Specifically, to create a list of protected words that would be guaranteed to remain intact, maintaining this way the core parts of the legal document that need to be static according to the authors.

Bibliography

- [AYS22] Yasuhiro Aoki, Masaharu Yoshioka, and Youta Suzuki. “Data-Augmentation Method for BERT-based Legal Textual Entailment Systems in COLIEE Statute Law Task”. In: *The Review of Socionetwork Strategies* 16.1 (2022), pp. 175–196.
- [CA17] Ilias Chalkidis and Ion Androutsopoulos. “A deep learning approach to contract element extraction”. In: *Proceedings of the 30th International Conference on Legal Knowledge and Information Systems (JURIX 2017)*. Luxembourg City, Luxembourg, 2017.
- [Cha+19] Ilias Chalkidis, Emmanouil Fergadiotis, Prodromos Malakasiotis, and Ion Androutsopoulos. “Large-Scale Multi-Label Text Classification on EU Legislation”. In: *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. 2019, pp. 6314–6322.
- [Cha+20] Ilias Chalkidis, Manos Fergadiotis, Prodromos Malakasiotis, Nikolaos Aletras, and Ion Androutsopoulos. “LEGAL-BERT: The Muppets straight out of Law School”. In: *Findings of the Association for Computational Linguistics: EMNLP 2020*. Online, 2020, pp. 2898–2904.
- [Cha+21] Ilias Chalkidis, Manos Fergadiotis, Dimitrios Tsarapatsanis, et al. “Paragraph-level Rationale Extraction through Regularization: A case study on European Court of Human Rights Cases”. In: *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 2021, pp. 226–241.
- [Cha+22] Ilias Chalkidis, Abhik Jana, Dirk Hartung, et al. *LexGLUE: A Benchmark Dataset for Legal Language Understanding in English*. 2022. CoRR: abs/2110.00976.
- [CK19] Ilias Chalkidis and Dimitrios Kampas. “Deep learning in law: early adaptation and legal word embeddings trained on large corpora”. In: *Artificial Intelligence and Law* 27.2 (2019), pp. 171–198.
- [CO21] Gergely Csányi and Tamás Orosz. “Comparison of data augmentation methods for legal document classification”. In: *Acta Technica Jaurinensis* 15.1 (2021), pp. 15–21.
- [Dai+23] Haixing Dai, Zhengliang Liu, Wenxiong Liao, et al. *AugGPT: Leveraging ChatGPT for Text Data Augmentation*. 2023. CoRR: abs/2302.13007.

- [Dev+18] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding”. In: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Minneapolis, Minnesota, 2018, pp. 4171–4186.
- [Eva+18] Richard Evans, David Saxton, David Amos, Pushmeet Kohli, and Edward Grefenstette. *Can Neural Networks Understand Logical Entailment?* 2018. CoRR: abs/1802.08535.
- [Fel05] Christiane Fellbaum. “Wordnet and Wordnets”. In: *Encyclopedia of Language and Linguistics*. Elsevier, 2005, pp. 2–665.
- [Fen+21] Steven Y. Feng, Varun Gangal, Jason Wei, et al. “A Survey of Data Augmentation Approaches for NLP”. In: *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. Online, 2021, pp. 968–988.
- [Kob18] Sосuke Kobayashi. “Contextual Augmentation: Data Augmentation by Words with Paradigmatic Relations”. In: *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*. New Orleans, Louisiana, 2018, pp. 452–457.
- [Lan+20] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, et al. *ALBERT: A Lite BERT for Self-supervised Learning of Language Representations*. 2020. CoRR: abs/1909.11942.
- [Lee+19] Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, et al. “BioBERT: a pre-trained biomedical language representation model for biomedical text mining”. In: *Bioinformatics* 36.4 (2019), pp. 1234–1240.
- [Lew+20] Patrick Lewis, Myle Ott, Jingfei Du, and Veselin Stoyanov. “Pretrained Language Models for Biomedical and Clinical Tasks: Understanding and Extending the State-of-the-Art”. In: *Proceedings of the 3rd Clinical Natural Language Processing Workshop*. Online, 2020, pp. 146–157.
- [Mik+13] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. *Efficient Estimation of Word Representations in Vector Space*. Scottsdale, Arizona, USA, 2013.
- [Min+20] Junghyun Min, R. Thomas McCoy, Dipanjan Das, Emily Pitler, and Tal Linzen. “Syntactic Data Augmentation Increases Robustness to Inference Heuristics”. In: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 2020, pp. 2339–2352.
- [PMA22] Dimitris Pappas, Prodromos Malakasiotis, and Ion Androutsopoulos. “Data Augmentation for Biomedical Factoid Question Answering”. In: *Proceedings of the 21st Workshop on Biomedical Language Processing*. Dublin, Ireland, 2022, pp. 63–81.
- [PMB13] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. *On the difficulty of training Recurrent Neural Networks*. 2013. CoRR: abs/1211.5063.
- [Raf+23] Colin Raffel, Noam Shazeer, Adam Roberts, et al. *Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer*. 2023. CoRR: abs/1910.10683.

- [Raj+16] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. “SQuAD: 100,000+ Questions for Machine Comprehension of Text”. In: *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*. Austin, Texas, 2016, pp. 2383–2392.
- [RJL18] Pranav Rajpurkar, Robin Jia, and Percy Liang. “Know What You Don’t Know: Unanswerable Questions for SQuAD”. In: *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Melbourne, Australia, 2018, pp. 784–789.
- [San+20] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. *DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter*. 2020. CoRR: abs/1910.01108.
- [SBG23] T. Y. S. S. Santosh, Philipp Bock, and Matthias Grabmair. *Joint Span Segmentation and Rhetorical Role Labeling with Data Augmentation for Legal Documents*. 2023. CoRR: abs/2302.06448.
- [SHB16] Rico Sennrich, Barry Haddow, and Alexandra Birch. “Neural Machine Translation of Rare Words with Subword Units”. In: *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Berlin, Germany: Association for Computational Linguistics, 2016, pp. 1715–1725.
- [SJS06] Marina Sokolova, Nathalie Japkowicz, and Stan Szpakowicz. “Beyond Accuracy, F-Score and ROC: A Family of Discriminant Measures for Performance Evaluation”. In: vol. Vol. 4304. 2006, pp. 1015–1021.
- [SKF21] Connor Shorten, Taghi M. Khoshgoftaar, and Borko Furht. “Text Data Augmentation for Deep Learning”. In: *Journal of Big Data* 8.1 (2021), p. 101.
- [Spa+20] Harold J. Spaeth, Lee Epstein, Jeffrey A. Segal, et al. *Supreme Court Database, Version 2020 Release 01*. Washington University Law. 2020.
- [Tsa+15] George Tsatsaronis, Georgios Balikas, Prodromos Malakasiotis, et al. “An overview of the BIOASQ large-scale biomedical semantic indexing and question answering competition”. In: *BMC Bioinformatics* 16.1 (2015), p. 138.
- [WZ19] Jason W. Wei and Kai Zou. “EDA: Easy Data Augmentation Techniques for Boosting Performance on Text Classification Tasks”. In: (2019). eprint: abs/1901.11196.
- [Yan+19] Ge Yan, Yu Li, Shu Zhang, and Zhenyu Chen. “Data Augmentation for Deep Learning of Judgment Documents”. In: *International Conference on Intelligent Science and Big Data Engineering*. Nanjing, China, 2019, pp. 232–242.
- [ZPT18] Henghui Zhu, Ioannis Ch. Paschalidis, and Amir Tahmasebi. *Clinical Concept Extraction with Contextual Word Embedding*. 2018. CoRR: abs/1810.10566.
- [ZZL16] Xiang Zhang, Junbo Zhao, and Yann LeCun. *Character-level Convolutional Networks for Text Classification*. 2016. CoRR: abs/1509.01626.

List of Figures

1.1	Image compression of Mordecai the snow dog, truncating the SVD at various ranks r and therefore generating artificial image data. Figure copied from the book Data-Driven Science and Engineering 2nd Edition, by Steven L. Brunton and J. Nathan Kutz	2
2.1	Word2vec augmentation example in a chunk of our SCOTUS legal texts. . . .	9
2.2	Two different cases of masked-lm predictions applied in a chunk of our SCOTUS legal texts.	11
2.3	Application of Back-Translation employing French, using a chunk of our SCOTUS legal texts. It is apparent that after back-translation has finished, a new artificial chunk has been generated. Also, the simplicity of the new chunk's vocabulary in comparison to the original chunk is noteworthy. . . .	12
2.4	Examples of generating sentences for data augmentation. Figure copied from Aoki et al. [AYS22]	16
3.1	A bar chart of the class distribution in SCOTUS	22
3.2	BERT input representation. The input embeddings are the sum of the token embeddings, the segmentation embeddings, and the position embeddings. In the final hidden state (when the second [CLS] is processed), the model captures the overall information from the entire sequence. Figure copied from Devlin et al. [Dev+18]	24
3.3	A bar chart class distributions in SCOTUS. The first figure shows the class distribution using the basic augmentation method, whilst the latter demonstrates the class distribution using "smart" augmentation.	31
3.4	Loss Learning Curves for the best and worse seed of each method. We can see that seeds within each method display similar trends, with word2vec having the largest margins. Overall, the model's performance seems to achieve better results in masked-LM and word2vec, achieving less overfitting. The loss learning curves are derived from the optimal model versions, specifically the ones from the 4th epoch counted in reverse. This selection is guided by the application of early stopping after 3 epochs, ensuring that subsequent epochs do not contribute to improvement.	33

3.5	Confusion Matrix # No Augmentation	34
3.6	Confusion Matrix # Standard Augmentation	35
3.7	Confusion Matrix # Smart Augmentation	35
3.8	F1 per class with standard augmentation, using the development results of our runs.	36
3.9	F1 per class with smart augmentation, using the development results of our runs.	37

List of Tables

- 3.1 LEGAL-BERT's performance using the original data for 5 different seeds. The last column "Results" takes the average of the 3 best-performing seeds. The best-performing seeds are determined by their results in the category of development macro-F1. The bold cells indicate the best-performing seeds and the total results. 28
- 3.2 Data Augmentation using Word2vec compared to the 3 baselines. The bold cells indicate the best-performing method in each category. Word2vec with the current settings (0.95 similarity threshold), appears to surpass the baselines in the test set but fails to do so in the development set. 28
- 3.3 Word2vec hyper-parameter tuning. The bold cells indicate the best-performing method in each category. All of these results were produced by seed #3 which was the best-performing seed. In this table, we compare different similarity thresholds as well as different amounts of generated artificial instances (as the second column involves tripling the data). Word2vec with 0.9 and 0.95 similarity thresholds appear to be the best options. "Dev" loss, corresponds to the loss of the development dataset. 29
- 3.4 A head-to-head comparison of Word2vec with different similarities. The generated results depict the average performance of all 5 seeds (following the same notion as in all the previously reported results). Word2vec with 0.9 similarity is the best method for the development data. The bold cells indicate the best-performing method in each category. 29
- 3.5 A head-to-head comparison of Back-Translation with different languages. The bold cells indicate the best-performing method in each category. Back-translation using only French as a pivot language beats back-translation that employed 3 languages (French included) and becomes the default back-translation for the remainder of our experiments. 30
- 3.6 A full comparison between all methods and baselines across all seeds. The bold cells indicate the best-performing method in each category. 30

3.7	A full comparison between all methods and baselines, across all seeds, using smart Augmentation. The bold cells indicate the best-performing method in each category. We observe that the results of the original data remain the best overall, whilst the margin between them and the 3 augmentation methods increases.	32
-----	--	----