

# Exploring the Promises of Transformer-Based LMs for the Representation of Normative Claims in the Legal Domain

**Reto Gubelmann**

University of St.Gallen  
Rosenbergstrasse 30  
9000 St.Gallen

**Peter Hongler**

University of St.Gallen  
Varnbuelstrasse 19  
9000 St.Gallen

**Siegfried Handschuh**

University of St.Gallen  
Rosenbergstrasse 30  
9000 St.Gallen

{reto.gubelmann,peter.hongler,siegfried.handschuh}@unisg.ch

## Abstract

In this article, we explore the potential of transformer-based language models (LMs) to correctly represent normative statements in the legal domain, taking tax law as our use case. In our experiment, we use a variety of LMs as bases for both word- and sentence-based clusterers that are then evaluated on a small, expert-compiled test-set, consisting of real-world samples from tax law research literature that can be clearly assigned to one of four normative theories. The results of the experiment show that clusterers based on sentence-BERT-embeddings deliver the most promising results. Based on this main experiment, we make first attempts at using the best performing models in a bootstrapping loop to build classifiers that map normative claims on one of these four normative theories.

## Contents

### 1 Introduction

Disagreements about normative claims are notoriously hard to resolve, and in some cases, they are even hard to recognize as such. For instance, consider (1). Do you think that a tax system that follows this principle is just?

- (1) It is just to tax people with the same income equally.

Example (1) illustrates what we mean by a normative claim: A moral judgment of some kind, that is, an assertion that something is either morally right or wrong. As we restrict our scope to tax law, the normative claims that we are interested in pertain to moral judgments of specific tax systems. Hence, while example (1) counts as a normative claim, example (2) does not count. While the latter is also about tax law, it does not make a claim about what

is just or unjust in this domain, but rather what is legal.

- (2) It is illegal not to pay one's taxes.

Furthermore, compare example (3). This example is not normative in the same explicit sense as (1): It does not directly make a claim about what is just; however, in contrast to (2), its standing directly depends on an explicitly normative claim such as (1). If one rejects the latter, one will reject (3) as well. We call statements of the kind of (1) *directly normative*, while we call statements of the kind of (3) *indirectly normative*.

- (3) We should attempt to create a tax system that taxes people with the same income equally.

In the discussion on international tax law, claims of the kind of (1) are regularly made, and even more often they figure implicitly in the arguments of legal scholars, say when statements of the kind of (3) are used without further questioning or without further argument. However, very often, the authors do not make explicit the normative perspective from which they are arguing (for instance, by stating and defending a claim of type (1) with explicit recourse to the political-philosophical literature in this domain). This is not the result of ill will. It is becoming increasingly difficult for students as well as for practitioners in the field to keep an overview on the different normative positions in the field. However, without such an overview, the debate threatens to lose sight of the central normative presuppositions of their debates. In the worst case, adherents of different normative positions will retreat into their normative bubbles and hence permanently hinder any truly rational debate about these topics.

To move towards improving this situation, we

explore the promises of using state-of-the-art LMs to cluster directly normative statements in tax law texts. More specifically, we explore the capacities of various transformer-based LMs to cluster directly normative statements that have been identified by an expert as belonging to one of these four views together. We use a variety of configurations for our clusterers, including different clustering algorithms with a range of parameters, and we test clustering based on specific words as well as on whole sentences. Furthermore, we explore the promises of using the models that have shown to allow for the best clusterers to initiate a bootstrapping loop to build effective classifiers of indirectly normative statements.

The task in focus of this article is both challenging and important. It is challenging because recognizing the normative background of a statement such as (1) requires expert knowledge, and even with such expert knowledge, genuine uncertainties remain in some cases. Furthermore, considered from a technical perspective, when compared to other clustering tasks, the amount of lexical overlap is substantially higher. If proponents of two different normative theories talk about tax justice, they mean something substantially different compared to a proponent of the procedural view. However, we are not confronted with clear-cut ambiguity, such as in the case of *bank*. Both mean to capture the same idea, they just understand that idea quite differently.<sup>1</sup>

The task is important because the subject matter that is addressed in such normative arguments is of central importance for democratic societies. What counts as a just taxation system directly influences the set-up of a taxation system and impacts the lives of the members of that society. Hence, providing support to navigate such normative landscapes is of central importance for societies – even more so if, as we have suggested, it is difficult for tax law researchers to be sensitive to these normative categories and the entire debate threatens to disintegrate into a number of normative bubbles.

After discussing related research, we introduce our datasets and experiment, we list the results and

---

<sup>1</sup>In the philosophical and linguistic debate, such concepts are called “essentially contested concepts”. The conception was first proposed by (Gallie, 1955), for recent discussions see (Collier et al., 2006) and (Rodriguez, 2015). According to this conception, concepts such as TAX JUSTICE are such that essential parts of their meaning are disputed. And the reason for the dispute is that the disagreement is due to larger-scale differences in worldview.

discuss them. We conclude by providing an outlook to further research.

## 2 Related Research

As mentioned above (section 1), we approach the clustering task in two different ways, one of them word-based, the other sentence-based. In the word-based case, our approach shares intriguing analogies to word-sense disambiguation (WSD), which is why we also introduce previous work in WSD. For the sentence-based approach, to the best of our knowledge, we cannot build upon more specific research than the generic and well-established clustering algorithms used for text and word clustering, which is why we focus on this research there. Furthermore, we also introduce the clustering algorithms with a focus on WSD, even though we will be using it to cluster not only words but also sentences – the simple reason being that it still seems to be a related task to cluster words that have different meanings and statements of different senses of central legal concepts such as tax justice.

### 2.1 Word-Sense Disambiguation

According to Navigli (2009, 3), WSD is the task to “computationally determine which sense of a word is activated by its use in a particular context”. Hence, in WSD, the goal is to decide, for any given word and context, which sense of the word in question is relevant for the context in question. Similarly, in our task of classifying normative statements, the goal is often to determine in what sense central concepts such as *tax justice* are meant in a given context.

Word-Sense Disambiguation (WSD) has been labelled an AI complete problem, that is, a problem that is as difficult to solve as general artificial intelligence as a whole. This claim, usually attributed to the seminal survey by Navigli (2009), but in fact first issued by Mallery (1988, 47), is as catchy as it is hard to cash out. The claim is in analogy to the concept of NP completeness (nondeterministic polynomial) from complexity theory, which refers to problems that can currently not be solved directly, but only using numeric approximations. Despite this vagueness, the label is accurate insofar as WSD is among the most notoriously difficult tasks in NLP, with systems still struggling to beat very primitive baselines such as the most frequent sense baseline, where an ambiguous term is simply always mapped onto its most frequent sense.

We expect the task of mapping words or sentences onto normative categories to be similarly challenging as the mapping of words to distinct senses.

**Supervised Approaches to WSD** In supervised approaches to WSD, the system learns a classifier based on large amounts of sense-annotated texts. An intuitive example of such a supervised approach is k-Nearest Neighbor (kNN) classification. In this approach, the features of the labelled data are stored as individual data points in the feature space. Any occurrence of a word to be disambiguated is then matched against all of the examples stored; the one sense having most examples within the class of k nearest data points in the feature space is selected as the correct sense in the specific context.

While there is a general consensus in the field that supervised systems perform best, they do so only within domains where there is sufficient high-quality labelled training data. Hence, recent years have seen semi-supervised approaches receiving more attention. In semi-supervised approaches, researchers typically try to interpolate the existing labels to unlabeled samples, for instance by assigning the sample the one label whose features are closest (in cosine distance) to the unlabelled sample in question.

Papandrea et al. (2017) present a toolkit for supervised word-sense disambiguation. They have developed a Java API that represents the WSD tasks as consisting of four subtasks – parsing, pre-processing, feature extraction, and classification. Researchers can replace the default elements that fulfill these subtasks with their own modules and train as well as test the resulting systems efficiently via the command line.

Duarte et al. (2021) present an in-depth analysis of semi-supervised approaches. Their general framework consists of four steps: (1) they extract pos-tags and word embeddings both from the target word to be disambiguated and from a fixed number of context words and combine them all into one single feature vector. (2) They create a graph out of all the feature vectors using a kNN approach. (3) They propagate the labels to the unlabeled data. (4) They use the resulting model in disambiguation challenges. What makes the paper stand out is that, within this overall framework, they vary a large number of parameters including the label propagation algorithm used - they employ the popular local and global consistency Zhou et al. (2004), label propagation Yamaguchi et al. (2016), Gaus-

sian fields Zhu et al. (2003) and OMNI-prop (Yamaguchi et al., 2015); they achieve the best results with label propagation. They also vary the embeddings used, with contextualized word embeddings from BERT (Devlin et al., 2019) and ELECTRA (Clark et al., 2020) performing best. Generally, the variation within a given benchmark is rather small, typically around 0.1 F1. Sousa et al. (2020) use a very similar approach, also testing different combinations of label propagation algorithms with embeddings.

**Knowledge-Based Approaches** In knowledge-based approaches to WSD, the occurrences of ambiguous words are disambiguated by relying on external knowledge. An early example of such a knowledge-based approach is the so-called *Lesk* algorithm (see (Lesk, 1986)). The basic idea is to measure the word overlap between the immediate context of the occurrence and the different definitions of the senses. The algorithm then selects the one sense where the overlap is largest.

Knowledge-based approaches, while usually not reaching the precision of supervised ones, have the advantage that they can rely on existing knowledge-bases such as WordNet (Miller, 1995). This makes knowledge-based approaches generally cheaper to extend to many different domains, as researchers do not have to create costly, time-consuming labelled datasets.

On a very general level, most of the systems follow the same structure (see (Scarlini et al., 2020), (Pasini et al., 2020), (Bevilacqua and Navigli, 2020) for recent work that follows this structure). First, the words whose meanings are to be disambiguated are identified in a dataset, and the CWE of their occurrences in the corpus are extracted. For each word to be disambiguated, the extracted CWE are then clustered, say, using a k-means algorithm. Then, using a knowledge-base such as WordNet (Miller, 1995) and a matching algorithm such as the page-rank based UKB (Agirre et al., 2014), the clusters are associated with the senses annotated in the knowledge base. In sum, this yields a dictionary of polysemous words, where each of the senses of each word is assigned a vector that characterizes the specific sense contextually. During runtime, the occurrence of the ambiguous words are then disambiguated by matching the CWE of the occurrence in question with the vectors stored and associated with each word: the occurrence of the word is mapped onto the one sense whose em-

bedding is closest to the CWE in question.

## 2.2 Clustering

From the perspective of WSD, clustering approaches are a kind of unsupervised approaches. In their purest forms, these clustering approaches are no typical cases of word sense disambiguation, but rather a cases of word sense discrimination: typical clustering approaches that function without tagged training data yield clustered representations of occurrences that allow to identify semantically related uses of words, without (in the absence of explicitly listed senses, say from a knowledge base) mapping these clusters onto specific senses.

A historical example of such an unsupervised word sense clustering algorithm is given by word spaces (see (Schütze, 1998) and again (Navigli, 2009, 26-28)). The underlying assumption of this approach is that ambiguities can be solved by looking at the co-occurrence patterns of ambiguous words. For instance, the biological sense of bug is supposed to co-occur with words that signal its biological constitution (food, habitat, reproduction, etc.), whereas the software sense is supposed to co-occur with IT vocabulary. Following this basic idea, each word is assigned a vector that represents its typical co-occurrences within a specified window (say, a number  $n$  words before and after the occurrence) in a given corpus, a so-called centroid, or context vector. The context vectors of an ambiguous word can then be clustered, say by agglomerative clustering, where, starting with singletons, iteratively further nearest members are added to the cluster until a threshold is reached.

A popular alternative to agglomerative clustering is k-means-clustering proposed by Lloyd (1982), where the centroids are initialized randomly and then iteratively updated. See (Géron, 2019, 238) for an overview and a non-technical recipe. The basic algorithm runs as follows:

1. Set the number of clusters  $k$  as a hyperparameter
2. Initialize the centroids randomly, say by randomly picking  $k$  instances.
3. Label the instances: Assign them to the closest centroid.
4. Update the centroids
5. ...

Over time, various improvements in efficiency to this simple algorithm have been made, such as trying to choose centroids that are far away from each other to initialize the algorithm (Arthur and Vassilvitskii, 2006), to avoid any unnecessary distance calculations (Elkan, 2003), and to process large datasets in batches (Sculley, 2010).

While k-means, especially when optimized as sketched, is very fast and scales well, it has its drawbacks. Notably, you need to specify the number of clusters  $k$  as a hyperparameter, the algorithm performs by design poorly on datasets that have clusters of various sizes and non-spherical shape (it uses a simple distance metric, after all).

Another popular clustering algorithm is DBSCAN (Ester et al., 1996). It's basic idea is to find regions of high density. One specifies as hyperparameters the maximal distance  $\epsilon$  and the minimal number of instances to be located within this distance. Then the algorithm finds core instances that have the minimal number of instances within  $\epsilon$ ; if one of these instances in turn counts as a core instance by having at least the minimal number of instances within  $\epsilon$ , it is added to the cluster formed by the initial core instance, and so on.

Advantages of DBSCAN are that it copes well with clusters of very different shapes, and it has just two hyper-parameters – no need to specify the number of clusters in advance, as with k-means.

In our main experiment, we are using k-means and DBSCAN with a range of different parameters, in our exploration of two classifiers to initiate a bootstrapping loop, we are using kNN-classifiers. In the future, we plan to also explore the promises of label-propagation in the sense of unsupervised approaches to WSD as well as the creation of a domain-specific knowledge-base that could then be used for knowledge-based approaches.

## 2.3 Transformer-Based LMs & Classical Word Embeddings

We use three different kind of model to deliver the embeddings for our experiments, which are all based on previous research: Transformer-Based LMs that deliver word-like embeddings, transformer-based LMs that have been fine-tuned with the specific goal to deliver high-quality sentence embeddings, and classical, non-transformer-based word embeddings. Since its publication in Vaswani et al. (2017), the transformer architecture has been very influential in virtually all domains

of NLP, including natural language understanding (NLU). With the exception of two models, all of the models tested are derived from this basic architecture.

**Word-based LMs** We here use three well-researched transformer-based LMs, namely bert-base-cased and bert-large-cased (Devlin et al., 2019) as well as roberta-large (Liu et al., 2019),<sup>2</sup>

**Sentence-based LMs** We test a large number of SBERT-Models (Reimers and Gurevych, 2019), as initial explorations showed that they perform clearly best. These SBERT-Models are based on a variety of transformer-based LMs (in addition to the classical BERT and RoBERTa, these are mpnet (Song et al., 2020), distilroberta (Sanh et al., 2019), xlm (Lample and Conneau, 2019), AIBERT (Lan et al., 2019), and minilm (Wang et al., 2020)).

**Classical WE** The classical, static word embeddings, namely GloVe (Pennington et al., 2014) and Komninos (Komninos and Manandhar, 2016), are included for purposes of comparison.<sup>3</sup>

### 3 Datasets

While we report the datasets for the main experiment and for the bootstrapping exploration separately, we here introduce the four normative categories that we have asked an expert to identify and that form the basis of both the clustering as well as the classifying experiment.

According to the so-called *Deontological View*, a tax policy proposal is just if it focuses on the treatment of the tax payer and not on the distribution of the income within a society. Hence, according to the Deontological View, one should neither look at democratic procedures, nor at the effects that a given tax system would have on the economy. Rather, one should look at whether it conforms to basic moral principles, such as the equality of all human beings. In this sense, example (1) constitutes a clear instance of the Deontological View.

According to the *Rawlsian View*, a tax system is just if it would be chosen by individuals that are under Rawls' famous veil of ignorance. Under this veil, individuals do not know their educational, financial, social, or any other position in the society whose tax system they are supposed to develop. It

is generally agreed that such individuals would favor tax systems strongly focused on equality – as they might end up as financially and educationally disadvantaged members of this society. For the purpose of the present analysis, a Rawlsian view emphasizes the need for redistribution from the rich to the poor. Of course, this is an oversimplification of Rawls theory.

Taxation should mainly result from good, democratically grounded processes – this is the gist of the *Procedural View*. Such view includes positions that argue for a certain tax policy proposal based on a discussion or debate about the arguments against and in favor of such proposal.

The fourth and final normative theory used in this article is the *Libertarian View*. According to it, taxation should be kept at a minimum in general, as it is considered illegitimate in all but a few cases. Obviously, this view strongly contrasts with the Rawlsian View, as the latter is much friendlier to redistribution of wealth, if it can be expected to contribute to equality.

#### 3.1 Main Experiment

For the clustering, we asked the expert to manually choose 10 samples of each of the four normative categories identified, trying to find directly normative rather than indirectly normative statements – where, of course, the boundary between the two is not always razor-sharp. This yields a total of 40 samples that were submitted to the clustering experiment. The samples are all grammatical sentences taken from publications in peer-reviewed journals from the legal domain, and their categorization was conducted by an expert. Then, a philosopher without special expertise in tax law went through the examples and annotated any disagreements. Where the disagreements could not be resolved by discussion, a different sample, whose categorization was uncontroversial, was chosen.

#### 3.2 Bootstrapping for Classifiers

The input for our kNN-Classifying experiment is given by four articles focusing on tax law, belonging to various text sorts, where the expert suspected – but did not antecedently identify – indirectly normative claims. **Article 1** is a tax-related discussion directed at the educated public,<sup>4</sup> **article 2** is a re-

<sup>2</sup>These models were downloaded from huggingface.co, see Wolf et al. (2019).

<sup>3</sup>The sbert- as well as the classical models were obtained from [https://www.sbert.net/docs/pretrained\\_models.html](https://www.sbert.net/docs/pretrained_models.html).

<sup>4</sup>See here: <https://www.americanprogress.org/issues/economy/reports/2020/09/28/490816/capital-gains-tax-preference-ended-not-expanded/>,

search article, akin to the articles from which the 40 samples were taken.<sup>5</sup> **Article 3** is a research article focusing on Chinese situations. This article has been published in a peer-reviewed journal focusing on tax issues in the Asia Pacific (Xu, 2021). **Article 4** is a memo of a parliamentary debate from Canada,<sup>6</sup> Given the samples used to create classifiers, we expect the classifiers to return the best results on article 2, then on article 3, then on article 1, and finally on the parliamentary memo, that is, on article 4.

#### 4 Experiment: Clustering Directly Normative Statements

For this clustering experiment, we took the 40 expert-chosen samples of directly normative statements and submitted it to a number of clustering algorithms. We tested different clustering algorithm (k-means vs. DBSCAN), each with a variety of parameters, different kinds of embeddings used for clustering, and a variety of pre-trained models to produce the embeddings. Details can be found on table 1.

**DBSCAN & k-means: parameters** As introduced above, the two clustering algorithms tested here are among the most commonly used for clustering. For k-means, one has to specify in advance the number of clusters that should be learned. As the target number of clusters is 4, we wanted to give the algorithm some wiggle room in case the four categories are in fact better mapped onto five clusters (say, because the embeddings of some model are more naturally clustered into five clusters). For DBSCAN, two parameters have to be set (see above, section 2), namely the maximal distance between two members of the same cluster (called epsilon), and the minimal number of members of one single cluster. As it is more difficult, in our case, to guess good values for these parameters, we have tested the algorithm with a much wider range of parameters. For epsilon, the range is 2, 2.5, ... 7, yielding 11 values, for the minimal number of members per cluster, the values tested are 2, 3 and 4.

**Models & Embedding Types** We are testing three different kinds of models; for references, see above, section 2.3; for the full list of models, see

last consulted on 13 August 2021.

<sup>5</sup>We selected section B (Kleinbard, 2016, 666ff.) for our testing.

<sup>6</sup><https://sencanada.ca/Content/SEN/Committee/362/bank/rep/rep003may00-e.htm>, last consulted on August 20, 2021.

thee appendix, table 3. We use four different routines to extract the embeddings:

**Word-Based** In this routine, we use a list of pre-compiled words that are intended to encapsulate central concepts of dispute, such as tax, taxation, or VAT (for the full list of these words, see the appendix). In this version of the experiment, the clusterers are not clustering sentences, but rather these words, as they appear in the sentences. As there is more input from experts, we expected these clusterers to perform better than the others. Here, we use well-researched transformer-based LMs, namely RoBERTa and BERT (see above, section 2.3)

**Sentence-Averaged Word-Based** In this routine, we use the average of all word embeddings, as the model delivers it for all words in the sentence. Hence, the sentence-embedding used here is the average of all word embeddings whose words appear in the sentence. Here, we use well-researched transformer-based LMs, namely RoBERTa and BERT (see above, section 2.3)

**Sentence-based** Here, we use the embeddings, as they directly result from the sentence-bert models trained by Reimers and Gurevych (2019). These models also output the average of all word embeddings (which we manually compute in the second variant), but they have been fine-tuned on the sentence level by training them on a wide variety of sentence-level tasks and datasets (the original models reported in (Reimers and Gurevych, 2019) use the combination of the SNLI and the Multi-Genre NLI datasets). Furthermore, the models that they fine-tuning are of many flavors, ranging from classical BERT to recent proposals such as mpnet (see above, section 2.3).

**Average of Classical Word Embeddings** We here test two classical kinds of word embeddings, GloVe as well as Komonos (see above, section 2.3), again taking the average of all word embeddings as the sentence embedding.

Overall, this resulted in 875 different clusterers. All of these clusterers were then given a simple task: to cluster the 40 sentences from the dataset.

algorithms	algorithm-specific parameters	model	embedding-type
DBSCAN	min. members (3), eps (11)	word-based transformers (3)	word-based
		sentence-bert (17)	mean of all words
		classical (2)	sentence-based
k-means	number of clusters (2)	word-based transformers (3)	word-based
		sentence-bert (17)	mean of all words
		classical (2)	sentence-based

Table 1: Overview on the configurations of the clustering systems tested. DBSCAN was used with three different minimal member counts, eleven different parameters for epsilon. k-means was tested with two different cluster sizes. Each of them was then tested with three different word-based LMs, which were each tested by focusing on a specific word as well as taking the mean of all words. Furthermore, we also tested 18 sentence-bert and two classical sentence-embedding models. Overall, this results in a total of 825 configurations for DBSCAN and 50 configurations for k-means tested. For details, see the appendix.

## 5 Results

We report the results of the clustering by giving what we call the average weighed homogeneity (AWH) of a given clusterer. The weighed homogeneity expresses the ratio of the largest member category and the overall membership in the cluster, weighed by the relative size of the cluster. For instance, if cluster 0 has 10 members, 8 of which are Ralwsian, then the homogeneity of the cluster would be 8/10, whereas the weighed homogeneity would be 8/10 multiplied by 10/40 (as the overall number of samples is 40), resulting in a weighed homogeneity score of 0.2. We take the average of the weighed homogeneity of all clusters produced by the clusterer. Without weighing the homogeneity, a classifier could cheat by performing well in very small clusters, but very poorly in one large cluster where most of the samples are.

Compare figure 1 for an overview on the results of this first experiment, showing the twenty best performing clusterers. For instance, the bar on the very left shows a score of 0.19, which represents the weighed average homogeneity that the model paraphrase-distilroberta-base-v2 could achieve using a k-means classifier (with specifications that are not shown on the chart, to facilitate overview). Overall, the chart shows that sbert-type embeddings perform clearly superior to both word and averaged-word embeddings. Surprisingly, embeddings not based on BERT, namely GLOVE-Average word embeddings land on rank 5 with a AWH-score of 0.17.

Figures 2 and 3 show the detailed clustering behavior of the two best performing clusterers. The former figure shows that the distilroberta-based embeddings allow for a surprisingly good clustering. Each of the four categories constitutes the largest

membergroup in one of the four clusters (the same holds for the multilingual-based clusterer, see figure 3). For instance, the first cluster listed in figure 2 contains 10 members in total, 9 of which belong to the procedural category. Furthermore, looking at figure 1, it is clear that these two models are not very exceptional, but rather the best of a number of similarly well-performing clusterers.

## 6 Discussion

### 6.1 Clustering

The results of the clustering experiment surpassed our optimistic expectations. The two best models deliver embeddings that allow for clusterers which obviously latch onto the characteristics of the four categories. Even without any explicit instruction, the clusterers identified the four categories, in some cases, as in the first cluster of distilroberta, with impressive precision.

Furthermore, experiment one shows that k-means outperforms DBSCAN in this task, a finding that was not affected by the very substantial variations in the parameters of DBSCAN used in the experiment. This is what one would expect given the data: As mentioned above (section 2.2), DBSCAN is good at finding areas of high densities. However, it seems that in this context, there are hardly any areas of high densities, and if there are, they do not latch onto the important categories. k-means, in contrast, is able to deal with distributions of various densities – given the right hyperparameter, it always delivers the corresponding number of clusters, even if the differences in densities within each cluster is very large.

Similarly, the sbert-models performed clearly best. Taking individual words or means over individual words, in contrast, performed poorer. On the

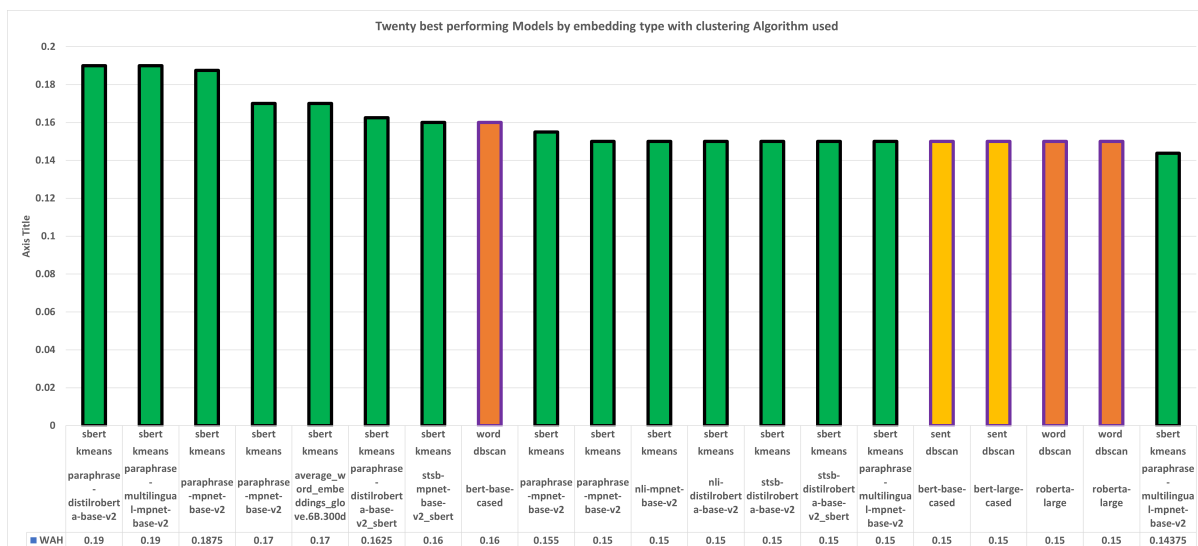


Figure 1: Results of the first experiment. Only the twenty best performing clusterers are shown. green filling represents sbert-type embeddings, red word-based embeddings, yellow word-based averaged embeddings. Black framing of bars represents k-means clustering, purple framing represents DBSCAN clustering

largest membergroup	total members	homog.	wh	lmg label
9	10	0.90	0.23	Procedural
8	12	0.67	0.20	Rawlsian
5	6	0.83	0.12	Deontological
8	12	0.67	0.20	Libertarianism

Figure 2: Detailed clustering behavior of distilroberta.

largest membergroup	total members	homog.	wh	lmg label
8	9	0.89	0.20	Rawlsian
9	13	0.69	0.23	Libertarianism
8	9	0.89	0.20	Procedural
6	9	0.67	0.15	Deontological

Figure 3: Detailed clustering behavior of multilingual.

one hand, this is plausible, as these models were developed specifically to deliver high-quality sentence embeddings. On the other hand, we expected the word-based approaches to perform well, too, as they were given further expert-compiled clues as to the specific words that are central for the task.

Finally, it is very interesting that a rather small model – one based on distilroberta – as well as a multilingual model outperform the large and monolingual models. This largely confirms the rankings on the sbert-page for clustering<sup>7</sup>, while it does not answer the question as to why smaller models are of better use for clustering algorithms than larger ones. The multilingual model, finally, invites multilingual explorations.

In short, what the results of the first experiment show is that the basic set-up of this approach is very promising.

<sup>7</sup>[https://www.sbert.net/docs/pretrained\\_models.html](https://www.sbert.net/docs/pretrained_models.html), last consulted on August 13, 2021.

## 6.2 Classification for Bootstrapping

Encouraged by the results of our main experiment, we decided to explore the promises of using the insights gained by the experiment to build classifiers that could initiate a bootstrapping loop between experts and the classifiers. Hence, the question that we are trying to answer with this final exploration is whether classifiers built on merely 40 samples could be used to mine legal texts and yield useful suggestions for further normative statements that could then, after having been reviewed by experts, be used to build better classifiers, etc.

As a consequence of the results obtained in the clustering tests, we decided to focus on two sbert-models. Using these models to obtain the embeddings, we then trained two simple kNN-classifiers (k=3) that classify a given sentence as belonging to one of the four categories according to the three nearest neighbors of that sentence’s embedding. To avoid classifying clearly non-normative sentences, we computed the centroid of the forty samples and empirically determined a threshold to automatically exclude any sentences that are beyond this threshold: Any embedding whose cosine similarity with the centroid is less than 0.6 is considered non-normative and not classified.

We then ran the dataset through the five articles mentioned and manually determined precision and recall. As this determination again requires expert knowledge in political philosophy, we asked an expert in the field of normative discussions in tax



law as well as a philosopher to independently tag the results.

The results are shown in table 2. We only counted as true positive a result if it also returned the correct categorization – in addition to simply correctly realizing that a sentence was normative.

**Article 1** Overall, the classifier has split the article up in 55 sentence (or sentence-like) elements. Of these, 3 are normative in either the Rawlsian, Procedural, Deontological, or Libertarian sense in focus here. The distilroberta-based classifier has returned 3 positives, two of which are true positives. The multilingual-based classifier has returned four positives, three of which are true positives

**Article 2** Overall, the classifier has identified 100 different sentences, 24 of which contain normative statements. The distilroberta-based classifier returns 14 positives, eight of which are true positives. The multilingual-based classifier returns 22 results, 11 of which are true positives.

**Article 3** Overall, the classifier has split the article up in 99 sentences. Of these, 6 are normative in the specific ways in focus here. The multilingual-based classifier has 14 positives, 5 of which are true positives, and 9 are false positives. The distilroberta-based models shows only 2 positives, none of which are true positives.

**Article 4** Overall, the classifiers have identified 100 sentences. Of these, 3 are normative in the specific ways in focus here. The distilroberta-based classifier has 2 positives, none of which are true positives. The multilingual-based classifier has 10 positives, two of which are true positives.

Classifier	Art.1	Art.2	Art. 3	Art. 4
distilroberta	0.67/0.67	0.57/0.3	0/0	0/0
multilingual	0.75/1	0.5/0.46	0.36/0.83	0.2/0.67

Table 2: Results of the classifying experiment. We report precision/recall.

The results of the two classifiers are encouraging, given the goal to build a classifier that can initiate a bootstrapping loop. In particular, it is notable that the two classifiers both deliver very few false positives, given the fact that in all four texts, normative claims were a small minority (roughly 10% per document). These figures seem well enough to initiate the bootstrapping loop envisaged above.

Furthermore, given the very simple calculation of the boundary between normative and non-

normative, this is further evidence that the embeddings used in these experiments are promising for further analyses of normative judgments: Simple geometric properties of them can be used to draw a good distinction between normative and non-normative.

## 7 Conclusion & Outlook

In this article, we have explored the promises of using well-known clustering and classifying approaches together with state-of-the-art transformer-based LMs to process normative statements in the legal domain. Our results indicate that this approach does indeed hold substantial promise.

As our next steps, we plan the following:

1. Using a bootstrapping loop, let experts and classifiers compile a large dataset for classifying.
2. Develop more sophisticated classifiers, systematically search the hyperparameter space.
3. Use adversarial attack strategies to discern whether the classifiers are latching on merely to shallow lexical cues, or whether they are actually building on more sophisticated representations.

## References

- Eneko Agirre, Oier Lopez de Lacalle, and Aitor Soroa. 2014. Random walks for knowledge-based word sense disambiguation. *Computational Linguistics*, 40(1):57–84.
- David Arthur and Sergei Vassilvitskii. 2006. k-means++: The advantages of careful seeding. Technical report, Stanford.
- Michele Bevilacqua and Roberto Navigli. 2020. Breaking through the 80% glass ceiling: Raising the state of the art in word sense disambiguation by incorporating knowledge graph information. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2854–2864.
- Kevin Clark, Minh-Thang Luong, Quoc V Le, and Christopher D Manning. 2020. Electra: Pre-training text encoders as discriminators rather than generators. *arXiv preprint arXiv:2003.10555*.
- David Collier, Fernando Daniel Hidalgo, and Andra Olivia Maciuceanu. 2006. Essentially contested concepts: Debates and applications. *Journal of political ideologies*, 11(3):211–246.

- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*.
- José Marcio Duarte, Samuel Sousa, Evangelos Milios, and Lilian Berton. 2021. Deep analysis of word sense disambiguation via semi-supervised learning and neural word representations. *Information Sciences*, 570:278–297.
- Charles Elkan. 2003. Using the triangle inequality to accelerate k-means. In *Proceedings of the 20th international conference on Machine Learning (ICML-03)*, pages 147–153.
- Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu. 1996. A density-based algorithm for discovering clusters in large spatial databases with noise. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD-96)*, pages 226–231.
- Walter Bryce Gallie. 1955. Essentially contested concepts. In *Proceedings of the Aristotelian society*, volume 56, pages 167–198.
- Aurélien Géron. 2019. *Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow: Concepts, tools, and techniques to build intelligent systems*. O’Reilly Media.
- Edward D Kleinbard. 2016. Capital taxation in an age of inequality. *S. Cal. L. Rev.*, 90:593.
- Alexandros Komninos and Suresh Manandhar. 2016. Dependency based embeddings for sentence classification tasks. In *Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics: human language technologies*, pages 1490–1500.
- Guillaume Lample and Alexis Conneau. 2019. Cross-lingual language model pretraining. *arXiv preprint arXiv:1901.07291*.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2019. Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942*.
- Michael Lesk. 1986. Automatic sense disambiguation using machine readable dictionaries: how to tell a pine cone from an ice cream cone. In *Proceedings of the 5th annual international conference on Systems documentation*, pages 24–26.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Stuart Lloyd. 1982. Least squares quantization in pcm. *IEEE transactions on information theory*, 28(2):129–137.
- John C. Mallery. 1988. Thinking about foreign policy: Finding an appropriate role for artificially intelligent computers. In *Master’s thesis, M.I.T. Political Science Department*.
- George A Miller. 1995. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11):39–41.
- Roberto Navigli. 2009. Word sense disambiguation: A survey. *ACM computing surveys (CSUR)*, 41(2):1–69.
- Simone Pappadrea, Alessandro Raganato, and Claudio Delli Bovi. 2017. Supwsd: A flexible toolkit for supervised word sense disambiguation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 103–108.
- Tommaso Pasini, Federico Scozzafava, and Bianca Scarlini. 2020. Clubert: A cluster-based approach for learning sense distributions in multiple languages. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4008–4018.
- Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543.
- Nils Reimers and Iryna Gurevych. 2019. Sentencebert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*.
- Philippe-André Rodriguez. 2015. Human dignity as an essentially contested concept. *Cambridge Review of International Affairs*, 28(4):743–756.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*.
- Bianca Scarlini, Tommaso Pasini, and Roberto Navigli. 2020. With more contexts comes better performance: Contextualized sense embeddings for all-round word sense disambiguation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3528–3539.
- Hinrich Schütze. 1998. Automatic word sense discrimination. *Computational linguistics*, 24(1):97–123.
- David Sculley. 2010. Web-scale k-means clustering. In *Proceedings of the 19th International Conference on World Wide Web*, pages 1177–1178.
- Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2020. Mpnet: Masked and permuted pre-training for language understanding. *arXiv preprint arXiv:2004.09297*.

Samuel Sousa, Evangelos Milios, and Lilian Berton. 2020. Word sense disambiguation: an evaluation study of semi-supervised approaches with word embeddings. In *2020 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *NIPS*.

Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. 2020. Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. *arXiv preprint arXiv:2002.10957*.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2019. Huggingface’s transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*.

Diheng Xu. 2021. Evaluation of value added tax exemption for small and low-profit enterprises in china—an analysis based on the principle of proportionality. *Asia-Pacific Tax Bulletin*, 27(3):1–7.

Yuto Yamaguchi, Christos Faloutsos, and Hiroyuki Kitagawa. 2015. Omni-prop: Seamless node classification on arbitrary label correlation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 29.

Yuto Yamaguchi, Christos Faloutsos, and Hiroyuki Kitagawa. 2016. Camlp: Confidence-aware modulated label propagation. In *Proceedings of the 2016 SIAM International Conference on Data Mining*, pages 513–521. SIAM.

Dengyong Zhou, Olivier Bousquet, Thomas N Lal, Jason Weston, and Bernhard Schölkopf. 2004. Learning with local and global consistency. In *Advances in neural information processing systems*, pages 321–328.

Xiaojin Zhu, John Lafferty, and Zoubin Ghahramani. 2003. *Semi-supervised learning: From Gaussian fields to Gaussian processes*. School of Computer Science, Carnegie Mellon University.

## A Details on Experiments

Table 3 shows the full name of the models used.

Word-Based Models SBERT-Models Classical Models
bert-base-cased bert-large-cased roberta-large
paraphrase-TinyBERT-L6-v2 paraphrase-distilroberta-base-v2 paraphrase-mpnet-base-v2 paraphrase-multilingual-mpnet-base-v2 paraphrase-MiniLM-L12-v2 paraphrase-MiniLM-L6-v2 paraphrase-albert-small-v2 paraphrase-multilingual-MiniLM-L12-v2 paraphrase-MiniLM-L3-v2 nli-mpnet-base-v2 nli-roberta-base-v2 nli-distilroberta-base-v2 distiluse-base-multilingual-cased-v1 stsb-mpnet-base-v2 stsb-distilroberta-base-v2 distiluse-base-multilingual-cased-v2 stsb-roberta-base-v2
average_word_embeddings_glove.6B.300d average_word_embeddings_komninos

Table 3: Overview on the models tested In clustering.

Following is a list of the words used for the word-based clustering routine, in this order (that is, words listed first take precedence):

1. taxation
2. tax\*
3. VAT
4. revenue-raising
5. redistribution
6. income
7. reward
8. pay