



UNIVERSITY OF AMSTERDAM

MSC ARTIFICIAL INTELLIGENCE
MASTER THESIS

Enriching Textual Data with Document Structure For Sentence Classification

by
MAARTEN DE JONGE
10002109

May 27, 2018

42 EC
01 July 2017, 31 December 2017

Supervisor:

Dr MAARTEN MARX

Assessor:

Dr TO BE DETERMINED



INFORMATION AND LANGUAGE PROCESSING SYSTEMS GROUP
INFORMATICS INSTITUTE

Contents

1	Introduction	2
2	Problem Statement	3
2.1	Dataset	6
2.2	Research Question	6
3	Related Work	9
4	Methodology	10
4.1	Unsupervised	11
4.1.1	Hierarchical Agglomerative Clustering	11
4.1.2	K-Means	12
4.2	Supervised	13
4.2.1	Convolutional Neural Networks	14
4.2.2	Difference between convolutional and recurrent neural networks	17
5	Experimental Setup	17
5.1	Dataset	18
5.2	Testing performance	18
5.3	CNN performance	18
5.4	Generalisation	19
6	Evaluation	19
6.1	Regularisation	19
6.1.1	Dropout	19
6.1.2	Weight decay	21
6.2	Models	23
6.3	Number of clusters	23
7	Conclusion	26

1 Introduction

A healthy democracy relies on a transparent political process that is open to the common populace

hier valt vast iets leuks te quoten van een politiek filosoof

. A common step towards achieving this is by publishing the proceedings of the parliament's meetings, such as the *Bundestag* in Germany¹ or the *Tweede Kamer* in the Netherlands². These proceedings can be useful in various ways, for example:

- Double-checking whether a certain politician's actions in the parliament are consistent with their public stance.
- Using them as a source of data for text analysis, such as classifying political ideology[1].
- Tracking the change in ideological leaning of a politician or party over time[2].

In each of these cases it would be hugely beneficial if the data was properly indexed. If you want to know John Doe's stance on immigration, you would ideally simply query a database for speeches by John Doe regarding the topic of immigration without having to manually skim over hundreds of documents to find the relevant speeches. It is unfortunate then that many of these documents are published solely in unstructured formats, such as PDF or plain text. Projects such as Political Mashup³ handle this by writing systems to parse and then index these documents. The semantic information required for indexing is currently recovered using rule-based methods. In the case of a PDF document, this is fairly challenging. The data often gets transcribed by a human typist, compiled to a PDF, and then goes back into a PDF decompiler for easier processing; this adds a lot of places where minor variations can occur in the output even though the document itself uses a consistent layout. Dealing with this in a rule-based system entails using either highly general rules that lead to a large probability of false positives, or a large amount of highly specific rules which can quickly lead to a spaghetti-like mess of special cases and is very fragile to unseen cases.

I propose that by using a small number of manually annotated documents as a dataset, a machine learning algorithm can learn to classify sentences in a way that allows it to segment a document into its constituent parts, while being more robust to noise than its rule-based counterpart. The common ways to do sentence classification (e.g. convolutional neural networks [3], recurrent neural networks or the simpler bag-of-words models) operate on sentences in a vacuum, considering

¹<https://www.bundestag.de/protokolle>

²<https://www.tweedekamer.nl/kamerstukken>

³<http://search.politicalmashup.nl/about.html>

only their linguistic contents and ignoring any contextual information that might be present. This is to be expected considering that most of the common datasets in this area really *are* just small bits of text in a vacuum; often-used datasets involve Twitter messages or short product reviews. In this case however, the sentences come from a document with a rich structure providing a lot of context. Anecdotally, as a human it is trivial to discern section headers in a document even when the document is in a foreign language; simply the fact that the section header might be printed in bold and centered rather than left-aligned gives it away. Incorporating this structural data into the learning process will hopefully increase the performance of the system, either by simply scoring better on the used metrics, or perhaps more indirectly by requiring less data or training time to achieve the same score.

2 Problem Statement

The German parliament, called the *Bundestag*, publishes the proceedings of their meetings, as an effort to open up the political process to the common people. These proceedings have been continuously published starting in 1949. Figure 1 shows a sample page from one of these proceedings; the left column contains a continuation of a speech from the previous page as well as two moderately sized speeches, while the right column contains a large number of very short speeches. Having a large corpus of political proceedings like this is wonderful and opens a lot of doors for research regarding political discourse. Figure 2 shows the same page seen in Figure 1, but with a number of regions of interest (manually) colored in. Unfortunately this data is only available as PDF files, which in terms of internal representation are entirely unstructured. That means that none of the rich structure highlighted in Figure 2 is actually present in a computer-usable way.

The central problem in this thesis is the extraction of speeches from the documents, transforming each document into a structured series of speeches that can be serve as a useful entry point into further research. As a first step, the structural layout (as in Figure 2) will be extracted using unsupervised clustering algorithms. The textual contents of the document are then fed into supervised text classifier, where each piece of text is augmented with the corresponding layout information previously obtained. More details on this process will be supplied in Section 4.

As annotating data by hand is an expensive process, a big focus is on limiting the amount of required training data as much as possible; it would be preferable if the system was able to learn sufficiently from a handful (say, less than 5) of hand-annotated files.

Präsident Dr. Norbert Lammert

- (A) Ich darf bereits jetzt darauf aufmerksam machen, dass ich nach Schließen des Wahlgangs die Sitzung für die Auszählung der Stimmen unterbrechen werde. Stellen Sie sich bitte darauf ein, dass das etwa eine Stunde dauern kann, weil ja ein doch relativ komplexer Wahlgang ausgezählt werden muss.

Ich eröffne die Wahl.

Liebe Kolleginnen und Kollegen, darf ich fragen, ob jemand im Saal ist, der seine Stimme noch nicht abgegeben hat? Oder hat jemand einen gesehen, den er dann nicht mehr gesehen hat und der seine Stimme noch abgeben könnte? – Dann schließe ich diesen Wahlgang und unterbreche die Sitzung bis zur Bekanntgabe des Ergebnisses der Wahl. Wir werden den Wiederbeginn der Sitzung rechtzeitig durch entsprechende akustische und optische Signale in den Immobilien des Bundestages ankündigen. Stellen Sie sich bitte darauf ein, dass es etwa eine Stunde dauern kann, bis wir diesen ja doch umfangreichen Wahlgang mit der gebotenen Sorgfalt ausgezählt haben.

Die Sitzung ist unterbrochen.

(Unterbrechung von 13.42 bis 14.52 Uhr)

Präsident Dr. Norbert Lammert:

Die unterbrochene Sitzung ist wieder eröffnet.

- (B) Liebe Kolleginnen und Kollegen, ich kann Ihnen das Ergebnis der Wahl der Stellvertreterinnen und Stellvertreter des Präsidenten bekannt geben: abgegebene Stimmkarten 626. Alle abgegebenen Stimmen waren gültig.

Von den abgegebenen Stimmen sind entfallen auf Peter Hintze 449 Jastimmen, 122 Neinstimmen und 51 Enthaltungen. In diesem Falle, was mich ein bisschen überrascht, waren 4 Stimmen ungültig. Das heißt, es gibt keine Stimmkarte, die insgesamt ungültig war, was ja doch auf eine gewisse Pfriffigkeit der neuen wie der alten Kollegen schließen lässt, aber bei einzelnen Wahlgängen ist das offenkundig anders. Noch einmal: 449 Jastimmen, 122 Neinstimmen, 51 Enthaltungen. Ich darf das mit Ihrem Einverständnis gleich mit der Frage an die jeweiligen Kolleginnen und Kollegen verbinden, ob sie die Wahl annehmen. Ich darf den Kollegen Hintze, der damit die notwendige Mehrheit erkennbar erreicht hat, fragen, ob er die Wahl annimmt.

Peter Hintze (CDU/CSU):

Ich bedanke mich. Ich nehme die Wahl an.

(Beifall bei der CDU/CSU sowie bei Abgeordneten der SPD und des BÜNDNISSES 90/DIE GRÜNEN)

Auf den Kollegen Johannes Singhammer sind bei 6 ungültigen Stimmen 442 Jastimmen, 115 Neinstimmen und 63 Enthaltungen entfallen. Auch er hat damit die notwendige Mehrheit eindeutig und klar erreicht. Ich darf ihn fragen, ob er die Wahl annimmt.

Johannes Singhammer (CDU/CSU):

(C) Ich danke für den Vertrauensvorschuss und nehme die Wahl gerne an.

(Beifall im ganzen Hause)

Präsident Dr. Norbert Lammert:

Die Kollegin Edelgard Bulmahn hat bei wiederum 6 ungültigen Stimmen 534 Jastimmen erhalten.

(Beifall im ganzen Hause)

50 Kolleginnen und Kollegen haben mit Nein gestimmt, 36 haben sich der Stimme enthalten. Frau Bulmahn, ich darf auch Sie fragen, ob Sie die Wahl annehmen.

Edelgard Bulmahn (SPD):

Auch ich bedanke mich für das Vertrauen, und ich nehme die Wahl gerne an.

(Beifall im ganzen Hause)

Präsident Dr. Norbert Lammert:

Auf die vorgeschlagene Kandidatin Ulla Schmidt sind 520 Jastimmen entfallen.

(Beifall im ganzen Hause)

66 Kollegen oder Kolleginnen haben mit Nein gestimmt, 35 haben sich der Stimme enthalten. 5 Stimmen waren ungültig. Ich bin zuversichtlich, Frau Schmidt, dass Sie die Frage ähnlich beantworten wie die bisher angesprochenen Kolleginnen und Kollegen.

Ulla Schmidt (Aachen) (SPD):

(D) Herr Präsident, Sie haben wie meistens recht. Ich nehme die Wahl an und bedanke mich für das große Vertrauen. Danke schön!

(Beifall im ganzen Hause)

Präsident Dr. Norbert Lammert:

Auf Petra Pau sind 451 Jastimmen entfallen,

(Beifall im ganzen Hause)

bei 113 Neinstimmen und 45 Enthaltungen. 17 Stimmen waren in diesem Wahlvorgang ungültig. Ich darf Frau Pau fragen, ob sie die Wahl annimmt.

Petra Pau (DIE LINKE):

Ja, Herr Präsident, ich nehme die Wahl gern an, und, liebe Kolleginnen und Kollegen, ich freue mich auf die weitere Zusammenarbeit.

(Beifall im ganzen Hause)

Präsident Dr. Norbert Lammert:

Schließlich darf ich noch das Wahlergebnis für Claudia Roth bekannt geben. Bei 14 ungültigen Stimmen hat sie 415 Jastimmen erhalten. Es gab 128 Neinstimmen und 69 Enthaltungen. Sie ist damit gewählt.

(Beifall im ganzen Hause – Claudia Roth [Augsburg] [BÜNDNIS 90/DIE GRÜNEN]:

Figure 1 – A sample page from one of the Bundestag proceedings.

(A) **Präsident Dr. Norbert Lammert**

Ich darf bereits jetzt darauf aufmerksam machen, dass ich nach Schließen des Wahlgangs die Sitzung für die Auszählung der Stimmen unterbrechen werde. Stellen Sie sich bitte darauf ein, dass das etwa eine Stunde dauern kann, weil ja ein doch relativ komplexer Wahlgang ausgezählt werden muss.

Ich eröffne die Wahl.

Liebe Kolleginnen und Kollegen, darf ich fragen, ob jemand im Saal ist, der seine Stimme noch nicht abgegeben hat? Oder hat jemand einen gesehen, den er dann nicht mehr gesehen hat und der seine Stimme noch abgeben könnte? – Dann schließe ich diesen Wahlgang und unterbreche die Sitzung bis zur Bekanntgabe des Ergebnisses der Wahl. Wir werden den Wiederbeginn der Sitzung rechtzeitig durch entsprechende akustische und optische Signale in den Immobilien des Bundestages ankündigen. Stellen Sie sich bitte darauf ein, dass es etwa eine Stunde dauern kann, bis wir diesen ja doch umfangreichen Wahlgang mit der gebotenen Sorgfalt ausgezählt haben.

Die Sitzung ist unterbrochen.

(Unterbrechung von 13.42 bis 14.52 Uhr)

Präsident Dr. Norbert Lammert:

Die unterbrochene Sitzung ist wieder eröffnet.

Liebe Kolleginnen und Kollegen, ich kann Ihnen das Ergebnis der Wahl der Stellvertreterinnen und Stellvertreter des Präsidenten bekannt geben: abgegebene Stimmkarten 626. Alle abgegebenen Stimmen waren gültig.

Von den abgegebenen Stimmen sind entfallen auf Peter Hintze 449 Jastimmen, 122 Neinstimmen und 51 Enthaltungen. In diesem Falle, was mich ein bisschen überrascht, waren 4 Stimmen ungültig. Das heißt, es gibt keine Stimmkarte, die insgesamt ungültig war, was ja doch auf eine gewisse Pfiffigkeit der neuen wie der alten Kollegen schließen lässt, aber bei einzelnen Wahlgängen ist das offenkundig anders. Noch einmal: 449 Jastimmen, 122 Neinstimmen, 51 Enthaltungen. Ich darf das mit Ihrem Einverständnis gleich mit der Frage an die jeweiligen Kolleginnen und Kollegen verbinden, ob sie die Wahl annehmen. Ich darf den Kollegen Hintze, der damit die notwendige Mehrheit erkennbar erreicht hat, fragen, ob er die Wahl annimmt.

Peter Hintze (CDU/CSU):

Ich bedanke mich. Ich nehme die Wahl an.

(Beifall bei der CDU/CSU sowie bei Abgeordneten der SPD und des BÜNDNISSES 90/DIE GRÜNEN)

Auf den Kollegen Johannes Singhammer sind bei 6 ungültigen Stimmen 442 Jastimmen, 115 Neinstimmen und 63 Enthaltungen entfallen. Auch er hat damit die notwendige Mehrheit eindeutig und klar erreicht. Ich darf ihn fragen, ob er die Wahl annimmt.

(C) **Johannes Singhammer (CDU/CSU):**

Ich danke für den Vertrauensvorschuss und nehme die Wahl gerne an.

(Beifall im ganzen Hause)

Präsident Dr. Norbert Lammert:

Die Kollegin Edelgard Bulmahn hat bei wiederum 6 ungültigen Stimmen 534 Jastimmen erhalten.

(Beifall im ganzen Hause)

50 Kolleginnen und Kollegen haben mit Nein gestimmt, 36 haben sich der Stimme enthalten. Frau Bulmahn, ich darf auch Sie fragen, ob Sie die Wahl annehmen.

Edelgard Bulmahn (SPD):

Auch ich bedanke mich für das Vertrauen, und ich nehme die Wahl gerne an.

(Beifall im ganzen Hause)

Präsident Dr. Norbert Lammert:

Auf die vorgeschlagene Kandidatin Ulla Schmidt sind 520 Jastimmen entfallen.

(Beifall im ganzen Hause)

66 Kollegen oder Kolleginnen haben mit Nein gestimmt, 35 haben sich der Stimme enthalten. 5 Stimmen waren ungültig. Ich bin zuversichtlich, Frau Schmidt, dass Sie die Frage ähnlich beantworten wie die bisher angesprochenen Kolleginnen und Kollegen.

(D) **Ulla Schmidt (Aachen) (SPD):**

Herr Präsident, Sie haben wie meistens recht. Ich nehme die Wahl an und bedanke mich für das große Vertrauen. Danke schön!

(Beifall im ganzen Hause)

Präsident Dr. Norbert Lammert:

Auf Petra Pau sind 451 Jastimmen entfallen,

(Beifall im ganzen Hause)

bei 113 Neinstimmen und 45 Enthaltungen. 17 Stimmen waren in diesem Wahlvorgang ungültig. Ich darf Frau Pau fragen, ob sie die Wahl annimmt.

Petra Pau (DIE LINKE):

Ja, Herr Präsident, ich nehme die Wahl gern an, und, liebe Kolleginnen und Kollegen, ich freue mich auf die weitere Zusammenarbeit.

(Beifall im ganzen Hause)

Präsident Dr. Norbert Lammert:

Schließlich darf ich noch das Wahlergebnis für Claudia Roth bekannt geben. Bei 14 ungültigen Stimmen hat sie 415 Jastimmen erhalten. Es gab 128 Neinstimmen und 69 Enthaltungen. Sie ist damit gewählt.

(Beifall im ganzen Hause – Claudia Roth [Augsburg] [BÜNDNIS 90/DIE GRÜNEN])

Figure 2 – The same page as in Figure 1, hand-annotated with interesting regions that play a significant role in understanding the layout. The red regions are the headers that signify that someone is starting a speech, and contain information about who the speakers is. The blue regions are little interruptions (*Heiterkeits*), often signifying approval or displeasure (the regularly seen *Beifall* means applause). The green regions indicate plain blocks of text.

2.1 Dataset

PDF files are unfortunately rather difficult to work with; being a vector-based format, they have no internal concept of words or sentences. All that’s available are instructions for drawing a certain character at a certain position. This means that even something as seemingly trivial as obtaining the lines of text from a PDF requires some fairly involved logic and heuristics (for instance, one would think that simply taking all characters on a page with the same y coordinate would be sufficient until realizing that many documents have a layout with two columns of text). This is dealt with by using the `pdftohtml` script from the Poppler PDF rendering library⁴. This script converts a PDF file to an XML file containing logical lines of text along with the coordinates and size of the line. Figure 3 shows an example of a portion of a PDF file and the corresponding XML. The contents of the `text` elements forms the plaintext that is fed into the classification algorithm. Although the layout information contained in the properties (`top`, `left`, `width`, `height` and `font`) could be used for the layout clustering, this would effectively tie the performance of the clustering to the performance of `pdftohtml`’s line-extraction heuristics. Instead, a separate system is used (the Apache `pdfbox`⁵ library for Java). By using the bounding boxes of individual characters as the base for clustering, the troubles of line-extraction can be fully bypassed.

The dataset was obtained through a rule-based system as described in the introduction, which annotates each `<text>` element of the XML files with a boolean flag indicating whether said element starts a new speech. The system was run on documents from the 18th electoral period of the *Bundestag*, consisting of 211 documents dating from 2013 to 2017. Together these documents contain 43,252 `<text>` elements indicating the start of speeches (that is, positive training samples), and 2,602,793 other elements (negative samples). This adds to a total of 2,646,045 training samples, taking up 503 MiB. This is a rather lopsided distribution (there are roughly 60 negative samples for each positive sample), which will have to be accounted for by, for instance, using stratified sampling. Figure 4 shows the distribution of the number of positive samples per file, giving a guideline as to how many files would have to be annotated to reach a desired amount of positive samples.

2.2 Research Question

As the application of supervised text classification is fairly basic and well-studied, academically speaking the interesting portion of this problem is figuring out what

⁴<https://poppler.freedesktop.org/>

⁵<https://pdfbox.apache.org/>

Dr. Norbert Lammert (CDU/CSU):
Herr Alterspräsident, lieber Kollege Riesenhuber, ich
nehme die Wahl gerne an.

(Beifall im ganzen Hause – Abgeordnete aller
Fraktionen gratulieren dem Präsidenten)

(a) A portion of the source PDF.

```
<text top="122" left="125" width="143" height="16" font="3">
  <b>Dr. Norbert Lammert </b>
</text>
<text top="122" left="269" width="83" height="17" font="4">
  (CDU/CSU):
</text>
<text top="142" left="125" width="328" height="17" font="4">
  Herr Alterspräsident, lieber Kollege Riesenhuber, ich
</text>
<text top="158" left="108" width="156" height="17" font="4">
  nehme die Wahl gerne an.
</text>
<text top="186" left="141" width="278" height="17" font="4">
  (Beifall im ganzen Hause – Abgeordnete aller
</text>
<text top="203" left="158" width="242" height="17" font="4">
  Fraktionen gratulieren dem Präsidenten)
</text>
```

(b) XML created by running `pdftohtml`, corresponding to the PDF excerpt in Figure 3a. The contents of the `text` elements are used as inputs for the classification algorithm; the layout data contained in the properties is not used, as a separate software pipeline is used for the unsupervised clustering.

Figure 3 – A sample excerpt from a source PDF, along with its XML representation created by `pdftohtml`.

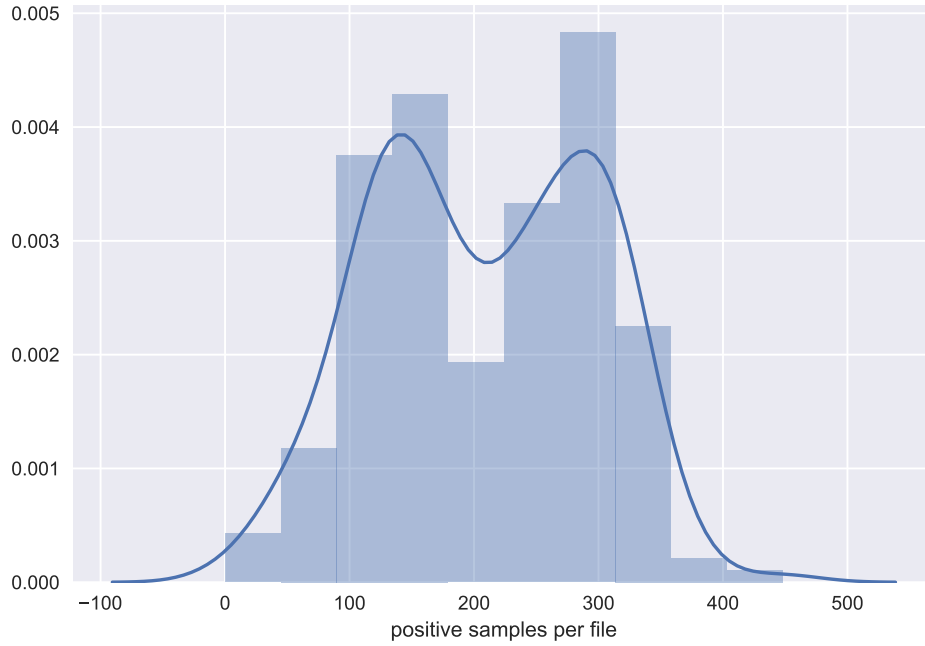


Figure 4 – Distribution of the number of positive samples per file

benefit is added by the unsupervised layout clustering, if any. A potential benefit could manifest itself in a couple of different ways:

- The performance (using a metric such as the F1 score or average precision, to be further elaborated on in Section 4) after training is increased.
- The peak performance is equal, but less training data is required to reach said peak.
- The peak performance is equal and a similar amount of training data is required, but the classifier needs less iterations to reach this performance.

This leads to the research question that is central to this thesis:

Research Question 1 *Does augmenting a text classification system with layout information obtained by unsupervised clustering of the input data improve the classifier with regards to:*

- *F1 score and average precision*
- *required volume of training data*

119 • *training speed*

120 3 Related Work

121 **Todo: Expand section**

122 In terms of analyzing document structure, Klampfl *et al.* [4] introduce a method
123 to analyze scientific articles, detecting blocks of text, labeling them (as e.g. section
124 headers, tables or references) and determining the reading order — all in an
125 unsupervised manner. The text block detection is done using a sequence of different
126 clustering algorithms, while the labeling is done using a heuristic approach.

127 The problem handled here is in a way similar to that of wrapper induction,
128 which is the process of inferring a *wrapper* (a program that extracts data into a
129 usable form) from a web page. A fairly recent survey of the state of this field is
130 done by Ferrara *et al.* [5], who note that a big problem is keeping up with the
131 constantly changing nature of web pages. A novel approach to combat this is that
132 of Gogar *et al.* [6]. They do wrapper induction by combining the textual content
133 of a webpage with a screenshot of the rendered webpage in an effort to do wrapper
134 induction on previously unseen web pages. The text is encoded in a way that
135 maintains spatial information, a model they refer to as *Text Maps* or *Spatial bag*
136 *of words*. The text maps and the screenshot are fed into separate convolutional
137 networks, after which the output is combined for a final classification. In a test of
138 extracting product names and prices from web pages, the system obtains very high
139 score comparable to systems that do use site-specific initialization.

140 Various forms of convolutional neural networks are commonly used for text
141 classification. The most basic architecture is described by Kim [3], where the input
142 words are tokenized and embedded before passing them to the convolutional neural
143 network. Additional exploration of the parameter space and its effect on various
144 datasets is done by Zhang & Wallace [7]. Comparable results are achieved by Zhang
145 *et al.* [8] by operating on the character-level rather than the word-level, bypassing
146 the issues overhead of using word embedding (either in extra training time or in
147 finding suitable pretrained embeddings). All the previously mentioned architectures
148 use a single convolutional layer; this is somewhat contrary to current trends in
149 computer vision, where popular models such as ResNet[9] go as deep as 152 layers.
150 This difference is explored by Conneau *et al.* [10], who take a character-level CNN
151 and show that adding more layers improves performance, before leveling out at 29
152 layers. They hypothesize that the difference in effective depth between computer
153 vision and language processing might be due to the difference in datasets. The
154 common ImageNet dataset used in computer vision deals with 1000 classes; in
155 contrast, sentiment analysis datasets vary between 2 and 25 classes. In addition,
156 they note that the deeper networks do require a larger amount of data to train.

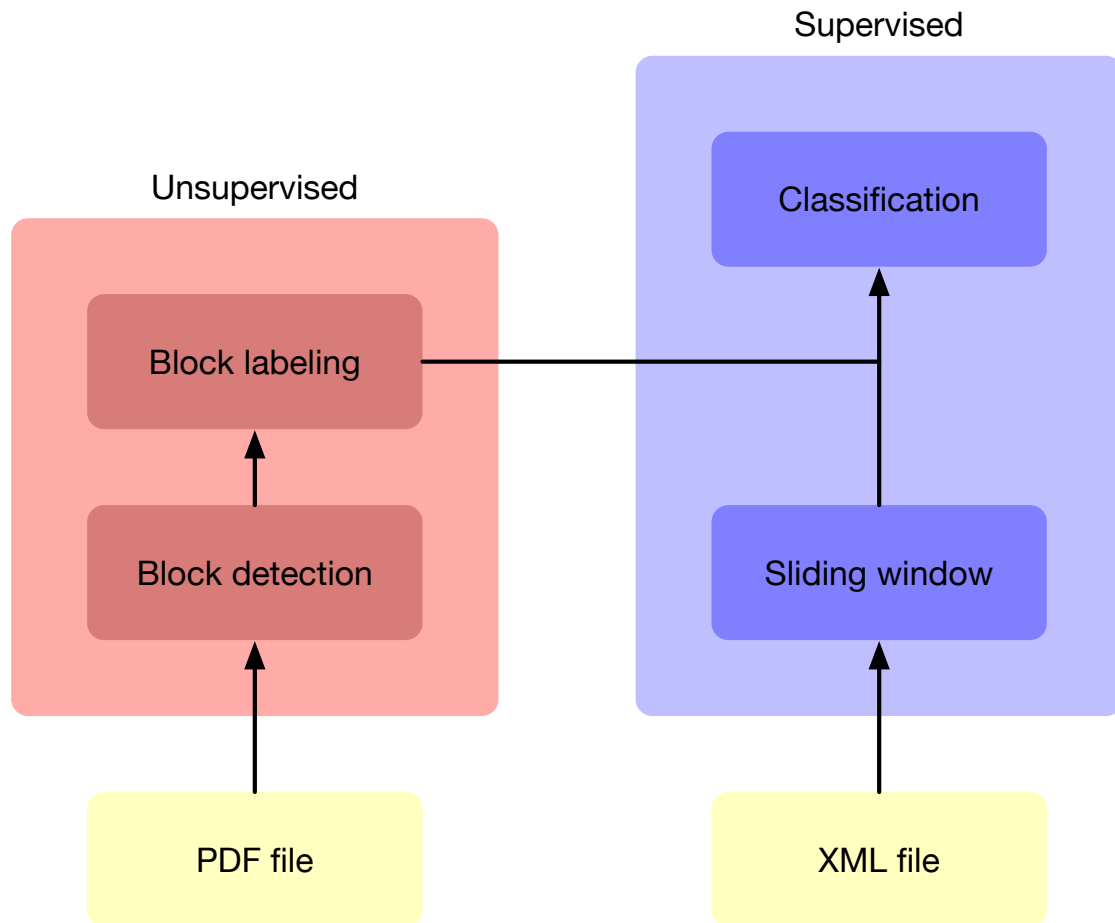


Figure 5 – A high-level overview of the system. The unsupervised block augments the input to the classifier.

157 4 Methodology

158 As described in section 2, the input data comes in the form of a PDF file which
 159 can be converted into plaintext XML data containing lines of text along with
 160 each line’s bounding box. There are two systems in play here; the main system
 161 is a convolutional neural network classifier that acts on the XML data. Further
 162 details on this are provided in section 4.2. The second system is the unsupervised
 163 preprocessor which acts on the PDF data and writes its output as an additional
 164 property into the XML data, which is elaborated upon in section 4.1. Figure 5
 165 shows a high-level overview of the two systems and how they interact.

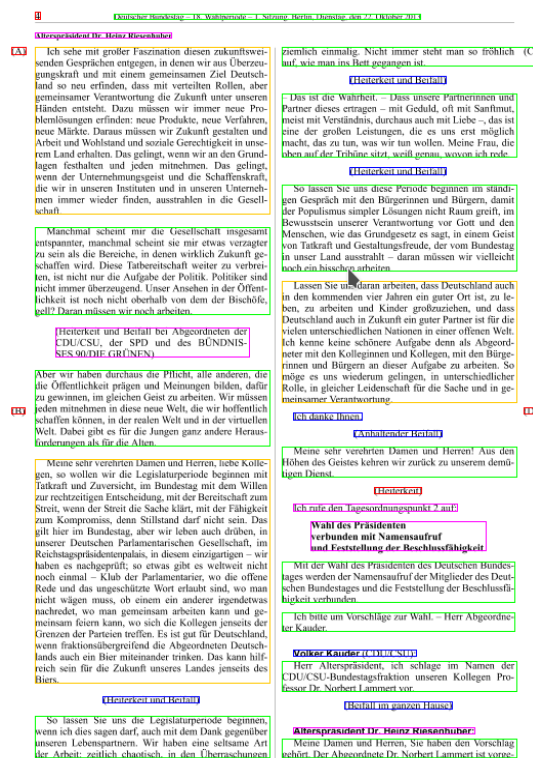


Figure 6 – An example of clustered blocks of text. Blocks with the same outline color belong to the same cluster.

166 4.1 Unsupervised

167 The unsupervised algorithm intends to detect and classify blocks of text in the
 168 PDF file; Figure 6 shows an example. This approach is based on work by Klampfl
 169 *et al.* [4], and consists of two separate clustering steps. First, individual characters
 170 (the fundamental objects available in a PDF file) are clustered together into blocks
 171 of semantically relevant text. These would be paragraphs, section headers, page
 172 decoration, etc. By using the bounding boxes of the blocks, they can be clustered
 173 based on their shape and some additional metadata (e.g. occurrence of font types
 174 and sizes). The following subsections will go into details on the two clustering
 175 steps.

176 4.1.1 Hierarchical Agglomerative Clustering

177 The first step is performed using hierarchical agglomerative clustering (HAC), an
 178 unsupervised bottom-up clustering algorithm that constructs a hierarchical tree
 179 of clusters (in this context referred to as a *dendrogram*). An example is shown in

180 fig. 7. The algorithm gets fed the individual characters present in the PDF files,
 181 then iteratively groups the two closest clusters (the initial inputs being regarded as
 182 clusters of one element) together until only a single cluster remains. This process
 183 involves two parameters:

- 184 1. The distance function between two characters.
- 185 2. The distance function between two clusters of characters.

186 The first parameter is trivially chosen to be the Euclidian distance between the
 187 coordinates of the two characters. The second parameter is called the *linkage* and
 188 has several common options, the most basic of which are:

- 189 • Single-linkage: The distance between clusters is based on the closest two
 190 elements:

$$d(A, B) = \min\{d(a, b) : a \in A, b \in B\}$$

- 191 • Maximum-linkage: The distance between clusters is based on the furthest
 192 two elements:

$$d(A, B) = \max\{d(a, b) : a \in A, b \in B\}$$

- 193 • Average-linkage: The distance between clusters is based on the average
 194 distance of its elements:

$$d(A, B) = \frac{1}{|A||B|} \sum_{a \in A} \sum_{b \in B} d(a, b)$$

195 As per Klampfl *et al.* [4], single-linkage clustering performs best for this task due
 196 to its tendency to form long thin clusters. This is beneficial since text is somewhat
 197 long and thin in nature (especially words and sentences). As an additional bonus,
 198 while the general time complexity for HAC is in $\mathcal{O}(n^3)$, single-linkage clustering
 199 can be done in $\mathcal{O}(n^2)$ [11], making it far more usable on realistic datasets.

200 After the dendrogram is constructed, it has to be cut at some level to obtain the
 201 desired blocks of text. Clustering can optionally be rerun using the newly found
 202 clusters as basic elements. This way, the document can incrementally be clustered
 203 from characters into words, words into lines, and finally lines into paragraphs.
 204 Both the level at which to cut the tree and the number of times to recluster are
 205 determined by trial and error based on the particular set of documents.

206 4.1.2 K-Means

207 The extracted blocks from the previous step are then clustered according to their
 208 similarity based on the following metrics:

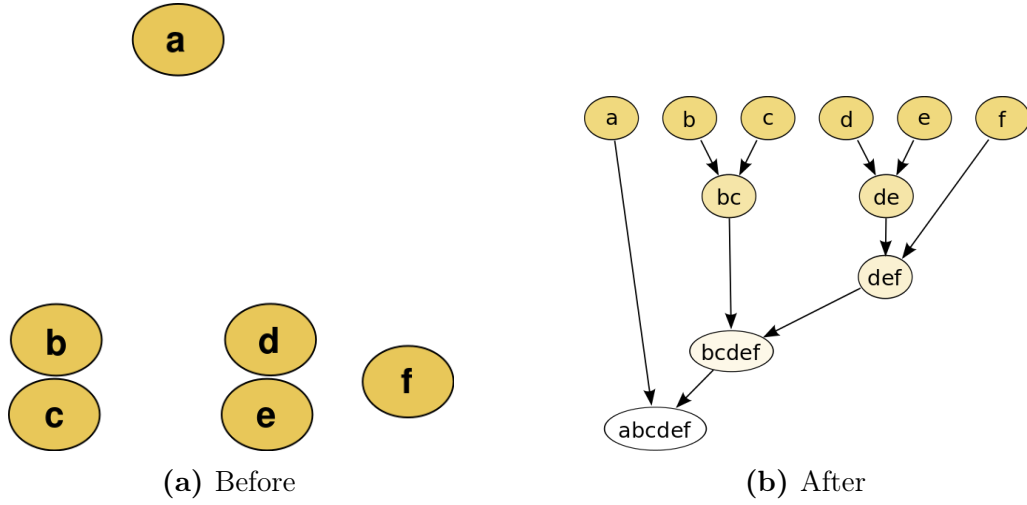


Figure 7 – An example of hierarical agglomerative clustering, where the nodes are clustered by distance.

- 209 • Width of the cluster
- 210 • Height of the cluster
- 211 • The ID of the most common font occurring in the cluster
- 212 • The size of the most common font occurring in the cluster

213 This is done using K-means clustering, with the value of k being varied for experi-
 214 mental purposes.

215 4.2 Supervised

216 After the data is augmented by the previously described clustering algorithms, it's
 217 fed into a convolutional neural network for classification. Since the source PDFs
 218 have a dual column layout with rather short lines, a sliding window is used to add
 219 valuable context, with the center element in the window supplying both the label
 220 (i.e. does or does not start a speech) and the cluster type. The text content of this
 221 window is then entered into a standard convolutional neural network architecture
 222 similar to the one proposed by Kim [3]. First an embedding layer is used to learn
 223 a high-dimensional representation of the words (no pre-trained embeddings were
 224 used because of both the specialized political domain of the data as well as a lack
 225 of quality pre-trained German models), followed by a number of Convolutional
 226 filters with 1-max pooling. The output of the filters is then concatenated into a
 227 single feature vector, which is amended with the cluster labels of each element in

the sliding window as an additional feature. This feature vector is then fed into a standard fully-connected neural network. This network features one hidden layer with a ReLU activation ($\text{relu}(x) = \max(0, x)$) and a single output node with a sigmoid activation. The layout of this system is detailed in fig. 8 along with its parameters and their baseline value.

The network is trained for 100 epochs, stopping early once no improvement in training loss has been made for 10 epochs. The optimisation process is done using the Adadelta[12] algorithm, with binary cross-entropy as the loss function. The loss will exhibit some up and down fluctuations after the global minimum has been reached; to account for this, the best loss so far along with the corresponding model parameters is kept track of throughout the learning process. Once all the epochs have been completed, the output is the model instantiated with the parameters corresponding to the lowest loss that was obtained. Regularisation is done through a combination of dropout and weight decay, the precise amounts to be determined through experimentation.

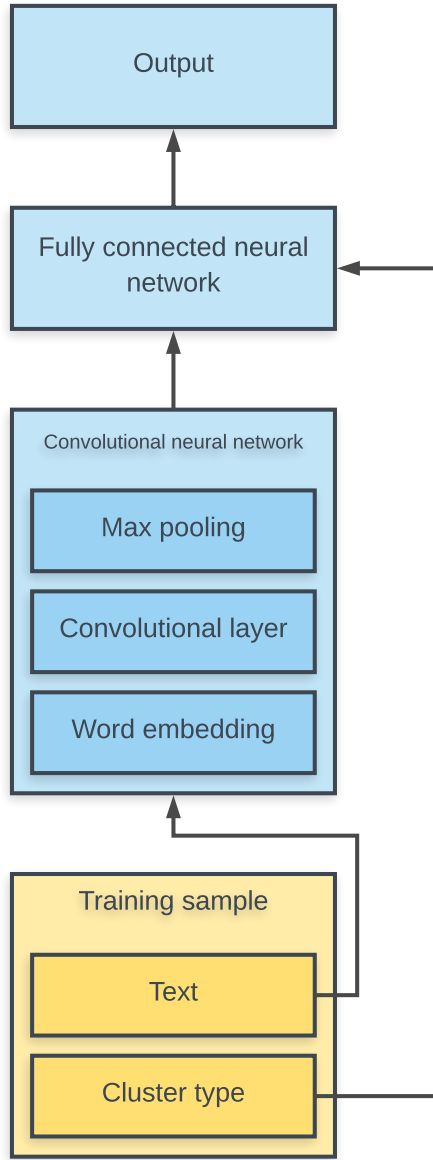
4.2.1 Convolutional Neural Networks

This is kind of jumbled and out of place, might need to be deleted or moved to an appendix or something.

For the purposes of machine learning, text produces 3-dimensional data: there is a feature vector for each word, and some (either variable or predetermined through padding) amount of words per text. This makes each training sample a 2-dimensional matrix, which then gets stacked in the depth dimension to produce 3-dimensional training data. This is troublesome as the standard machine learning algorithms work on 2-dimensional data, assuming a feature vector for each sample rather than a matrix. There are three common methods to deal with this:

- Bag of words
- Convolutional neural networks
- Recurrent neural networks

Aside from being completely different methods, they differ in a major way in how they handle the sequential nature of text. The bag of words approach is the simplest in that it simply disregards this sequential nature, instead creating what is essentially a histogram of word occurrences. This downsamples each sample from a feature matrix to a feature vector, allowing the use of normal machine learning algorithms (commonly support vector machines). While the sequential information can be kept to some degree by use histograms of n -grams rather than



(a) The layout of the supervised portion of the system. The input data contains text and a cluster type assigned by the unsupervised portion. The text gets put into a convolutional neural network, the output of which is fed together with the cluster type into a fully connected neural network. The sigmoid function is applied to the output of this final neural network to obtain the classification.

Parameter	Default value
Window size	9
Word embedding size	300
Number of filters	100
Filter sizes	3, 5, 7
Activation	ReLU
Pooling type	1-max

(b) The default parameters used in the convolutional neural network.

Parameter	Default value
Number of hidden layers	1
Hidden layer size	50
Number of outputs	1
Output activation	Sigmoid

(c) The default parameters used in the fully-connected neural network.

Figure 8 – The model and its parameters.

words (unigrams), this causes the size of the input data to scale exponentially with the value of n .

Convolutional neural networks (CNNs) work by taking a number of filters (sometimes called kernels or feature maps) of a specified size and convolving these over the input data. A simplified example using one filter is shown in fig. 9. In this

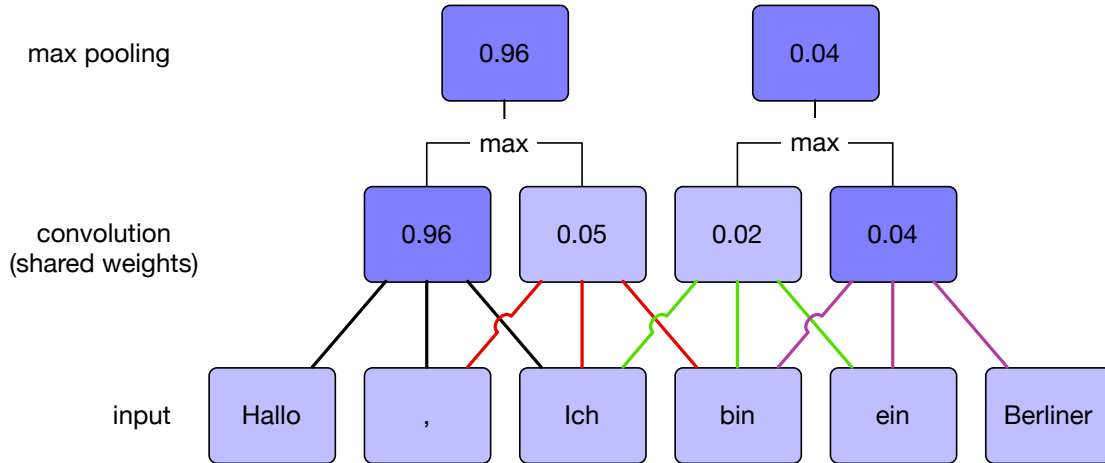


Figure 9 – A simplified convolutional neural network with one filter and max pooling.

example, the input text is convolved with a filter with a width of 3 and a stride of 1 — that is, each application of the filter considers three subsequent input elements, after which the window is shifted one space to the right. This filter is essentially a small neural network mapping three items to one output value, whose weights are reused for each application of the filter. Reusing the weights in this way (weight sharing) prevents the number of parameters in the network from spiraling out of control. [13] After the application of this convolution layer, the responses of the filter form a new sequence of roughly the same size as the input (minus a few due to missing the edges). The next step is to downsample this sequence by means of a *max pooling* layer, which maps all values within a window to the maximum value amongst those values. While conceptually similar to a convolution, this step generally does not involve overlap, instead either setting the stride to the same value as the window size (usually 2) reducing the entire sequence to 1 value (1-max pooling). The reason for this is twofold:

1. It downsamples the number of inputs, reducing the amount of parameters required further on in the network.
2. It adds translation invariance to the feature detected by this filter. The example filter of fig. 9 appears to react strongly to the presence of the word “Hallo”. Without the pooling layer, changing the location of the word

286 “Hallo” in the input would similarly change the location of the high activation
287 in the intermediate representation; this would be *equivariance*. The more
288 aggressively the pooling is applied, the higher the degree of invariance.

289 This combination of convolution followed by pooling can be repeated multiple times
290 as desired or until there is only a single value left as output from the filter. Finally,
291 the outputs of all filters are concatenated and fed into a standard feedforward
292 neural network.

293 While CNN architectures in computer vision are generally very deep, they
294 tend to be very shallow in natural language processing; commonly just a single
295 convolution followed by 1-max pooling [7].

296 4.2.2 Difference between convolutional and recurrent neural networks

297 Recurrent neural networks (in particular LSTMs or GRUs) are seemingly the most
298 natural fit for language processing, since they process an entire sequence and are
299 therefore fully conditioned on the word order (as opposed to the convolutional neural
300 networks which tend to learn translation invariant ngram features). Regardless,
301 convolutional neural networks will be used in this research. This choice is based on
302 two factors:

- 303 1. In practise, the performance for classification tasks does not differ between
304 the two types of networks.[14]
- 305 2. The computations in convolutional networks are highly independent of each
306 other, allowing for great paralellization (in particular with regards to running
307 on a GPU). In contrast, LSTMs are bottlenecked by the fact that each
308 calculation is dependent on the previous calculations. As a result, CNNs
309 achieve far higher training speeds.[15]

310 5 Experimental Setup

311 Experiments are focused on the difference in performance between the baseline
312 CNN model without clustering information (referred to from here on as **CNN**) and
313 the model augmented with clustering information (which will be referred to as
314 **CNN-cluster**). Performance will be measured with regards to the following three
315 metrics:

- 316 1. Number of training epochs until convergence
- 317 2. The F1 score or average precision metrics on a test set (see Section 5.2)

318 3. The number of training samples required to attain a specific F1/average
319 precision score

320 In each case, the experiment will be repeated 10 times by means of 10-fold cross
321 validation followed by a Student’s T-test to gauge the probability of the following
322 null hypothesis being true:

323 **Null Hypothesis 1** *Adding clustering information to the CNN model does not*
324 *change the performance of the model.*

325 5.1 Dataset

326 Referring back to Figure 4, the average document has between 100 and 300 positive
327 samples. Since a secondary concern is to minimize the number of documents that
328 would have to be annotated as training data, the tested dataset sizes will be very
329 low, with the number of positive samples being one of 100, 200, 500 and 1000. Due
330 to the relative abundance of negative samples and to prevent overfitting on the
331 distribution of the labels, stratified sampling will be used to keep a 1:1 ratio of
332 positive to negative samples. In addition to the size, the number of cluster types
333 (the k in k -means) will be varied to examine its impact on the performance.

334 5.2 Testing performance

335 All models will be tested on a test set containing 1000 positive samples and
336 10000 negative samples, all of which are guaranteed not to be in the training set.
337 Performance on this set is measured by constructing a precision-recall curve, and
338 calculating two values:

- 339 1. The average precision (which is equivalent to the area under the curve)
- 340 2. The F1 score of the point on the curve maximizing the F1 score

341 5.3 CNN performance

342 Although less central to the thesis than the difference between the CNN and
343 CNN-cluster models, some experimentation will be done with the parameters
344 of the convolutional network in an attempt to optimize the performance. These
345 parameters include the dimensionality of the word embeddings, the number of
346 filters, the pooling strategy (1-max versus a smaller region) and the number of
347 convolutional layers.

348 5.4 Generalisation

349 This particular dataset has the quirk that the performance of a rule-based system
350 created based on recent documents decreases in performance when used on older
351 documents, the older the document the worse it performs. This occurs despite the
352 layout being visually the same all the way back to the 1950s. A number of files
353 from old election periods has been labeled (and manually verified for correctness)
354 in order to test

- 355 1. whether the CNN models handle this better than the rule-based system does.
- 356 2. Whether the clustering-augmented CNN model performs better on this task
357 than the baseline CNN.

358 6 Evaluation

359 6.1 Regularisation

360 Since regularisation is very important to prevent overfitting, the common methods
361 of regularisation are tested first to obtain a good baseline value for the rest of the
362 experiments. These methods are:

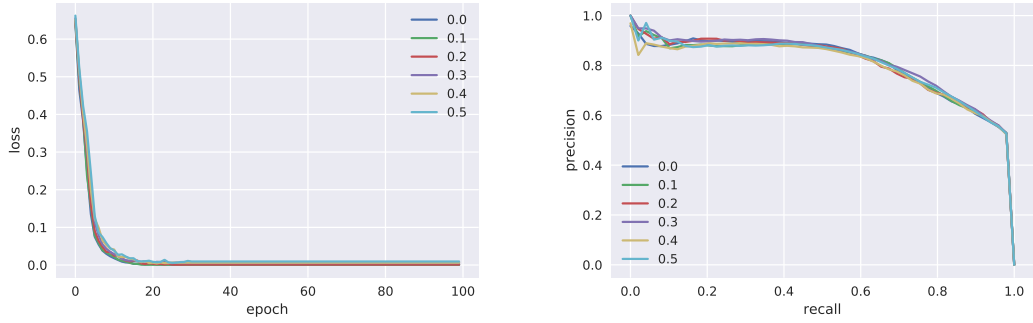
- 363 • Dropout, with values taken from the set $\{0.0, 0.1, 0.2, 0.3, 0.4, 0.5\}$.
- 364 • Weight decay, with values taken from the set

$$\{0.0, 1 \times 10^{-5}, 1 \times 10^{-4}, 1 \times 10^{-3}, 1 \times 10^{-2}\}.$$

365 Further details on these methods are given in section 4.

366 6.1.1 Dropout

367 Figure 10 shows the training loss over time, as well as the resulting precision-recall
368 curve on the training set. There does not seem to be any visible difference in
369 performance; this lack of effect is also shown in fig. 11. A dropout rate of 0.5
370 appears to score the best, although none of the distributions are different in a
371 statistically significant way. The full data, with the mean and standard deviation
372 of the F1 score and the area under the curve, are given in table 1.



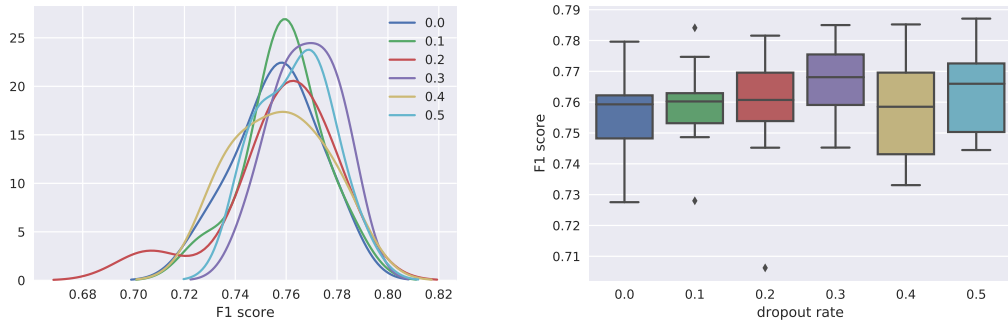
(a) The average loss at each epoch.

(b) The average precision-recall curves at various dropout values.

Figure 10 – The loss over time (a) and the precision-recall curve (b). In both cases, the values are averaged over 10 trials.

dropout rate	F1 mean	F1 stddev	AoC mean	AoC std	Area under averaged curve
0	0.756437	0.0150751	0.719548	0.00967498	0.806093
0.1	0.759067	0.0142657	0.718028	0.00749939	0.803368
0.2	0.757629	0.0199559	0.717309	0.0167461	0.807958
0.3	0.767517	0.012102	0.725084	0.0205084	0.815534
0.4	0.757488	0.0168819	0.715304	0.0118686	0.795153
0.5	0.76351	0.0131205	0.714167	0.0183001	0.804967

Table 1 – The F1 and AoC scores at various dropout values.



(a) Kernel density estimation

(b) Boxplot

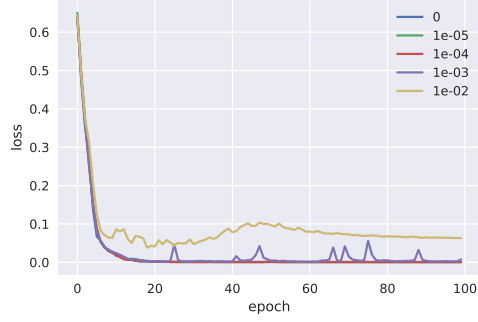
Figure 11 – A kernel density estimation and boxplot, based on the F1 score values over 10 repeated trials.

373 6.1.2 Weight decay

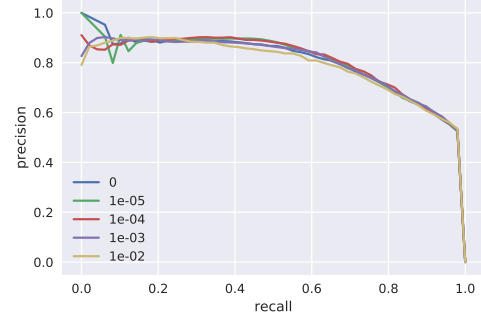
374 Similarly to the dropout plots, fig. 12 shows the training loss over time and the
 375 precision-recall curve at various weight decay values. There are again no significant
 376 differences, with a value of 1×10^{-4} performing just slightly better.

decay value	F1 mean	F1 stddev	AoC mean	AoC std	Area under averaged curve
0	0.759594	0.008119	0.710192	0.0078634	0.805901
1e-05	0.763445	0.0171368	0.715502	0.0113409	0.807963
0.0001	0.766299	0.0160939	0.732198	0.0102998	0.80583
0.001	0.762132	0.01537	0.736607	0.0131718	0.801772
0.01	0.753685	0.013023	0.737005	0.0150467	0.791353

Table 2 – The F1 and AoC scores at various decay values.

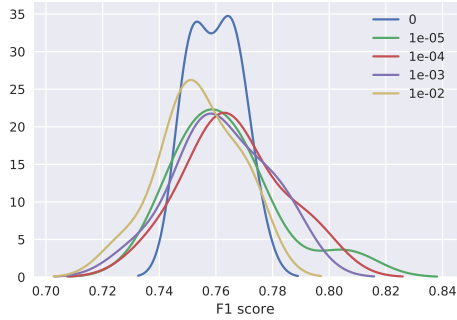


(a) The average loss at each epoch.

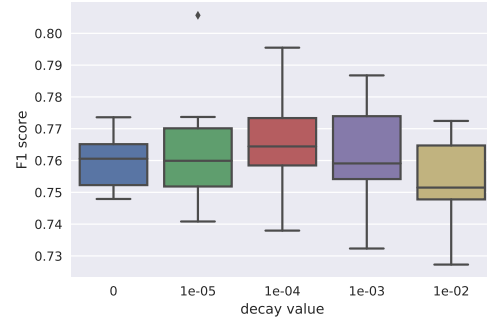


(b) The average precision-recall curves at various decay values.

Figure 12 – The loss over time (a) and the precision-recall curve (b). In both cases, the values are averaged over 10 trials.



(a) Kernel density estimation



(b) Boxplot

Figure 13 – A kernel density estimation and boxplot, based on the F1 score values over 10 repeated trials.

6.2 Models

Here the models described earlier (TODO: beschrijven en ernaar verwijzen) are tested with the default parameters of the CNN. Based on the previous results, the dropout rate was set to 0.5 and the decay to 1×10^{-4} . The plots in Figure 14 show a subtle but observable difference; the model using the full window jumps out of the pack in both cases, converging a little bit faster and being on top for most of the precision-recall curve, although the model with an additional CNN catches up at higher recalls.

The boxplot in fig. 15b paints a clear picture with three tiers of performance: the model without cluster labels performs worst, adding only the cluster of center of the window improves it a bit, while using the full window performs best with and without a CNN applied to it. Figure 15 shows an interesting relation between the two top performing models. Their score distributions are clearly bimodal with both peaks occurring at roughly the same scores, indicating the possibility of getting caught in a local optimum. The real win of the simpler model without the CNN appears to be that it is far less likely to get stuck in that local optimum.

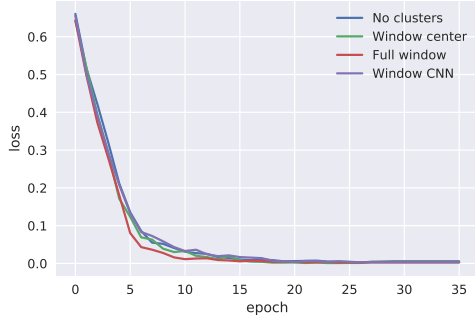
Significance values for the F1 score are given in table 4, which shows that every model is different from every other model using the common cutoff of $p \leq 0.05$. The least convincing is the difference between the two best models (Full window and Window CNN), but given the still significant result it is reasonable to state that the Full window model is the clearly winner given its far more reliable ability to reach the global optimum.

6.3 Number of clusters

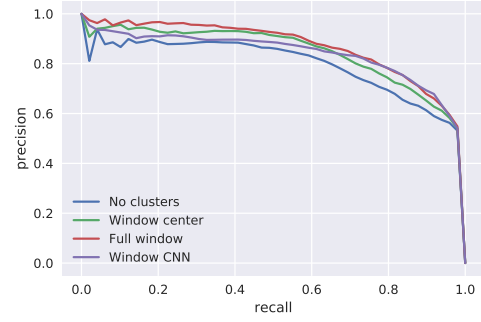
Figure 16 shows the results of varying the number of clusters created in the final k-means clustering step.

model	F1 mean	F1 stddev	AoC mean	AoC std	Area under averaged curve
No clusters	0.753191	0.0139397	0.715827	0.0125601	0.795299
Window center	0.785527	0.0081693	0.739108	0.0113586	0.843328
Full window	0.80462	0.0125957	0.758653	0.0210983	0.867003
Window CNN	0.804459	0.0120204	0.743269	0.0160533	0.838726

Table 3 – The F1 and AoC scores of the various models.



(a) The average loss at each epoch.

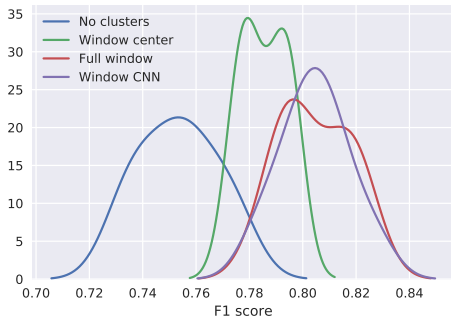


(b) The average precision-recall curves for the various models.

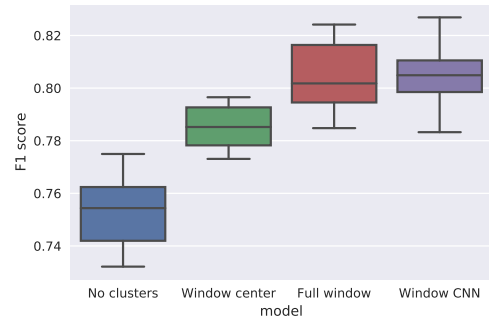
Figure 14 – The loss over time (a) and the precision-recall curve (b). In both cases, the values are averaged over 10 trials.

	No clusters	Window center	Full window	Window CNN
No clusters	nan	2.09153e-05	2.02693e-05	3.6803e-05
Window center	2.09153e-05	nan	0.00484301	0.00629707
Full window	2.02693e-05	0.00484301	nan	0.974039
Window CNN	3.6803e-05	0.00629707	0.974039	nan

Table 4 – The probability for each pair of models that their F1 scores are generated by the same underlying distribution (i.e. the probability of the null hypothesis that the models perform the same being true).



(a) Kernel density estimation



(b) Boxplot

Figure 15 – A kernel density estimation and boxplot, based on the F1 score values over 10 repeated trials.

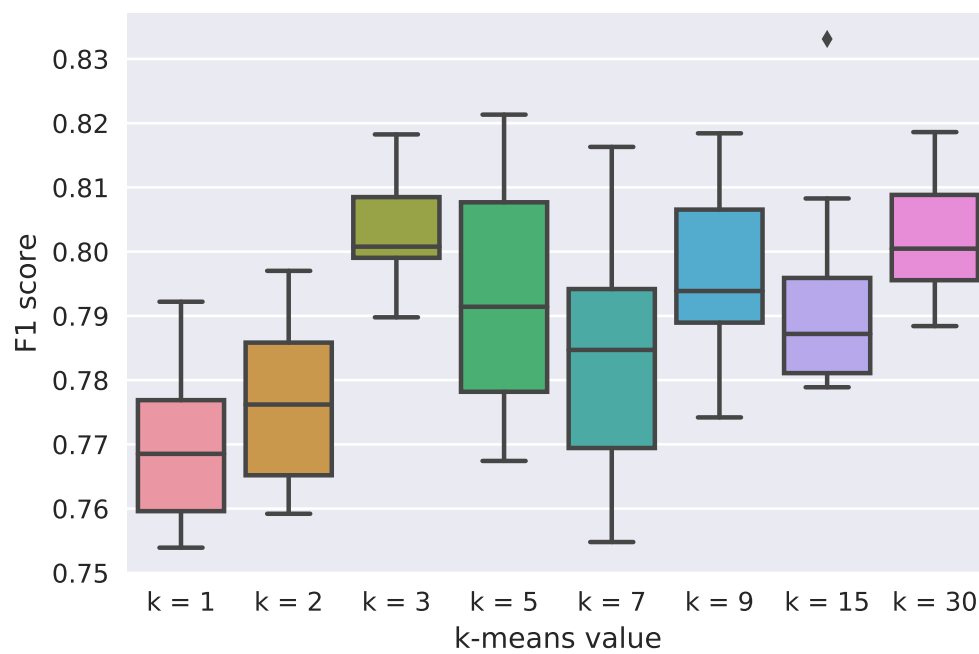


Figure 16 – A boxplot based on the F1 scores of 10 repeated trials for various amounts of clusters created by the k-means algorithm.

7 Conclusion

Todo: conclusion

References

1. Iyyer, M., Enns, P., Boyd-Graber, J. & Resnik, P. *Political ideology detection using recursive neural networks* in *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* **1** (2014), 1113–1122.
2. Gross, J. H., Acree, B., Sim, Y. & Smith, N. A. Testing the Etch-a-Sketch Hypothesis: A Computational Analysis of Mitt Romney’s Ideological Makeover During the 2012 Primary vs. General Elections (2013).
3. Kim, Y. Convolutional Neural Networks for Sentence Classification. *CoRR abs/1408.5882*. arXiv: 1408.5882. <<http://arxiv.org/abs/1408.5882>> (2014).
4. Klampff, S., Granitzer, M., Jack, K. & Kern, R. Unsupervised document structure analysis of digital scientific articles. *International Journal on Digital Libraries* **14**, 83–99 (2014).
5. Ferrara, E., Meo, P. D., Fiumara, G. & Baumgartner, R. Web Data Extraction, Applications and Techniques: A Survey. *CoRR abs/1207.0246*. arXiv: 1207.0246. <<http://arxiv.org/abs/1207.0246>> (2012).
6. Gogar, T., Hubacek, O. & Sedivy, J. *Deep Neural Networks for Web Page Information Extraction in Artificial Intelligence Applications and Innovations* (eds Iliadis, L. & Maglogiannis, I.) (Springer International Publishing, Cham, 2016), 154–163. ISBN: 978-3-319-44944-9.
7. Zhang, Y. & Wallace, B. C. A Sensitivity Analysis of (and Practitioners’ Guide to) Convolutional Neural Networks for Sentence Classification. *CoRR abs/1510.03820*. arXiv: 1510.03820. <<http://arxiv.org/abs/1510.03820>> (2015).
8. Zhang, X., Zhao, J. & LeCun, Y. *Character-level convolutional networks for text classification* in *Advances in neural information processing systems* (2015), 649–657.
9. He, K., Zhang, X., Ren, S. & Sun, J. Deep Residual Learning for Image Recognition. *CoRR abs/1512.03385*. arXiv: 1512.03385. <<http://arxiv.org/abs/1512.03385>> (2015).

- 435 10. Conneau, A., Schwenk, H., Barrault, L. & LeCun, Y. Very Deep Convolutional
436 Networks for Natural Language Processing. *CoRR* **abs/1606.01781**. arXiv:
437 1606.01781. <<http://arxiv.org/abs/1606.01781>> (2016).
- 438 11. Sibson, R. SLINK: an optimally efficient algorithm for the single-link cluster
439 method. *The computer journal* **16**, 30–34 (1973).
- 440 12. Zeiler, M. D. ADADELTA: An Adaptive Learning Rate Method. *CoRR*
441 **abs/1212.5701**. arXiv: 1212.5701. <<http://arxiv.org/abs/1212.5701>>
442 (2012).
- 443 13. LeCun, Y., Bengio, Y., *et al.* Convolutional networks for images, speech, and
444 time series. *The handbook of brain theory and neural networks* **3361**, 1995
445 (1995).
- 446 14. Yin, W., Kann, K., Yu, M. & Schütze, H. Comparative Study of CNN and
447 RNN for Natural Language Processing. *CoRR* **abs/1702.01923**. arXiv: 1702.
448 01923. <<http://arxiv.org/abs/1702.01923>> (2017).
- 449 15. Gehring, J., Auli, M., Grangier, D., Yarats, D. & Dauphin, Y. N. Convolutional
450 Sequence to Sequence Learning. *CoRR* **abs/1705.03122**. arXiv: 1705.03122.
451 <<http://arxiv.org/abs/1705.03122>> (2017).