



UNIVERSITY OF AMSTERDAM

MSC ARTIFICIAL INTELLIGENCE
MASTER THESIS

Enriching Textual Data with Document Structure For Text Classification

by
MAARTEN DE JONGE
10002109

October 12, 2018

42 EC
01 July 2017, 31 December 2017

Supervisor:

Dr MAARTEN MARX

Assessor:

Dr MAARTEN VAN SOMEREN



INFORMATION AND LANGUAGE PROCESSING SYSTEMS GROUP
INFORMATICS INSTITUTE

Contents

1	Introduction	2
2	Problem Statement	4
2.1	Dataset	7
2.2	Research Question	9
3	Related Work	9
4	Methodology	10
4.1	Unsupervised	10
4.1.1	Hierarchical Agglomerative Clustering	11
4.1.2	Clustering	13
4.1.3	Dirichlet Process Clustering	14
4.2	Supervised	15
4.2.1	Optimisation process	19
4.2.2	Regularisation	21
4.2.3	Convolutional Neural Networks	21
5	Experimental Setup	23
5.1	Dataset	23
5.2	Testing performance	23
5.3	CNN performance	24
5.4	Generalisation	24
6	Evaluation	24
6.1	Parameter Exploration	25
6.2	Training Set Size	26
7	Conclusion	29

Abstract

Proceedings of parliamentary debates are often published as unstructured PDF files, making them unsuitable for indexing into a database or querying for specific information. Digitizing them into a structured format is challenging; doing so manually is labor-intensive, doing so through a rule-based system is error-prone. Manually annotating a small number of documents can provide a training set that allows a convolutional neural network to accurately classify the elements that denote the structure of the document (e.g. the start of a new speech), using purely the textual content of the document. Adding a preprocessing step that clusters pieces of text based on the physical layout of the document improves the classification performance in a minor, but statistically significant way.

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 debates, 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 debates 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

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

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

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

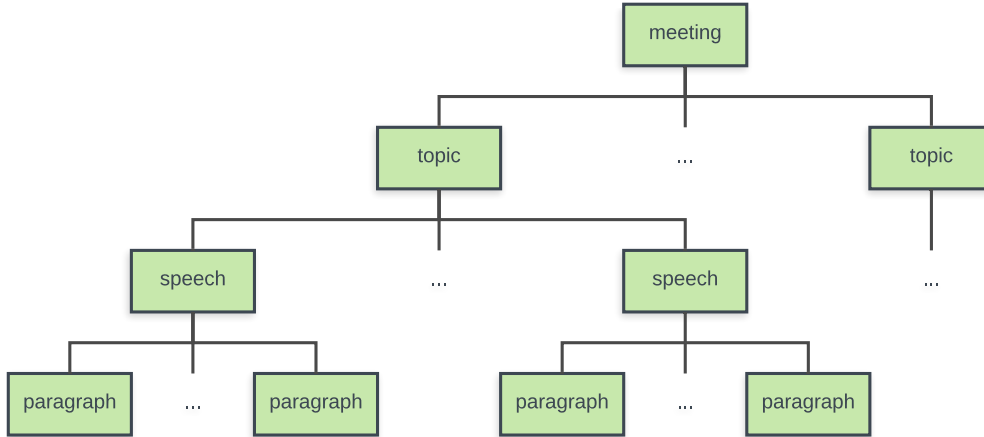


Figure 1 – The structure of a parliamentary debate is that of a shallow tree.

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 fragments of text 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. While the structure of a parliamentary debate is tree-shaped (a general example is shown in fig. 1), this tree is generally both shallow and very rigidly structured. As a result, simply classifying the positions in the text where a new structural element begins is enough to reconstruct the full tree given the domain knowledge about this particular set of documents; some variations can occur depending on the particular conventions used by the creator of the documents, for example in how speech interruptions by noisy colleagues are handled. After the document’s layout has been extracted, it is a matter of obtaining each element’s metadata (such as the speaker and their political affiliation for each speech element). This step is outside of the scope of this thesis.

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 only their linguistic contents

and ignoring any contextual information that might be contained in the layout in which the text might have been embedded. 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, using a clustering approach inspired by Klampfl *et al.* [4], 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 since the inception of the German Empire in 1871, when the parliament was called the *Reichstag*, up to 1942; publishing resumed after the conclusion of World War 2 with the inception of the Bundestag in 1949. Of course, the further back in the time you go, the more difficult the documents become to use. While the more recent documents are fully digital, documents prior to 1998 are scans of physical documents and require optical character recognition, which becomes even more problematic in the documents from the 1800s where a thick Gothic font is used. Figure 2 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. Figure 3 shows the same page seen in fig. 2, but with a number of regions of interest (manually) colored in. Unfortunately this data is only available in a PDF format, which in terms of internal representation is entirely unstructured. That means that none of the rich structure highlighted in fig. 3 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 fig. 3) 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

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 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.

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.

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 – 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 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):

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):

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])

(C)

(D)

Figure 3 – The same page as in Figure 2, 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 speaker 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.

103 this process will be supplied in section 4.

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

108 2.1 Dataset

109 The dataset contains 11 hand-annotated documents, comprising 2052 positive
110 samples and 76,944 negative samples. The average document contains about
111 200 positive samples. Seeing as the source documents are PDF files, some
112 transformations are needed to extract text from them in a way that is suitable
113 for use as a machine learning dataset. This is not as trivial as it might seem,
114 given that PDF files only contain instructions for drawing certain characters
115 at certain coordinates, with no internal concept of paragraphs, lines or even
116 words.

117 For the supervised portion of the system, the PDF files are transformed
118 using the `pdftohtml` utility from the Poppler PDF rendering library⁴. This
119 utility takes a PDF file and uses a number of heuristics to output lines of
120 text occurring inside the file. In this context, a *line* refers to any number of
121 characters that occur on the same height on a page while preserving reading
122 order (which is relevant when dealing with documents that have a two-column
123 layout) and is explicitly not the same as a sentence. These lines are output
124 in an XML format which includes metadata on the geometry of the line
125 (position and size) and its font. An example of a portion of text from the
126 dataset and the corresponding XML output from `pdftohtml` is shown in
127 fig. 4. Since this process is based on heuristics, it can and does go wrong;
128 there are instances of mistakes such as a single line being broken up into
129 multiple XML nodes, or lines of two side-by-side columns being taken as a
130 single XML node. This is not terribly common (the frequency of such errors
131 is perhaps one per file on average), but it does make the dataset inherently
132 noisy and is one of the issues that rule-based systems have trouble with.

133 For the unsupervised clustering, the PDF files are taken directly as input
134 and clustering is done using individual characters as the basic unit. This is
135 just as easy as clustering on the lines produced by `pdftohtml`, but has the
136 benefit of not being dependent on the imperfect line-extraction heuristics.

⁴<https://poppler.freedesktop.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">
  Dr. Norbert Lammert
</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 4a. 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 4 – A sample excerpt from a source PDF, along with its XML representation created by `pdftohtml`.

2.2 Research Question

A baseline model can be created by leaving out the clustering step, leaving us with two models:

1. A baseline model that classifies based on purely text
2. A model that classifies based on both text and layout information

There are two ways in which the second model can improve upon the baseline: either it performs better (using a metric such as the F1 score), or it requires less data to reach the same performance. This naturally leads to two research questions:

Research Question 1 *Does augmenting a text classification system with layout information obtained by unsupervised clustering of the input data improve the F1 score of the classifier?*

Research Question 2 *Does augmenting a text classification system with layout information obtained by unsupervised clustering of the input data allow the classifier to reach its peak performance using less input data?*

3 Related Work

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

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

For text classification, various forms of convolutional neural networks are commonly used. The most basic architecture is described by Kim [3], where

the input is tokenized and the tokens are embedded into a higher-dimensional representation before passing them to the convolutional neural network. Using this same architecture, Zhang & Wallace [7] perform additional exploration of the parameters and their effects on various datasets. Comparable results are achieved by Zhang *et al.* [8] by operating on the character-level rather than the word-level, bypassing the overhead of using word embedding (either in extra training time or in finding suitable pretrained embeddings) as well as freeing the researcher from having another layer in their network to tune. All the previously mentioned architectures use a single convolutional layer; this is contrary to current trends in computer vision, where popular models such as ResNet[9] go as deep as 152 layers. This difference is explored by Conneau *et al.* [10], who take a character-level CNN and show that adding more layers improves performance, before leveling out at 29 layers. They hypothesize that the difference in effective depth between computer vision and language processing might be due to the difference in datasets. The common ImageNet dataset used in computer vision deals with 1000 classes; in contrast, sentiment analysis datasets vary between 2 and 25 classes. In addition, they note that the deeper networks do require a larger amount of data to train.

4 Methodology

As described in section 2, the input data comes in the form of a series of PDF files. From this point on, there are two different systems acting on the data:

1. An unsupervised clustering step.
2. A supervised classifier.

The clustering step acts on the raw PDF file to segment it into blocks of text; these blocks are then clustered based on the shape of the blocks. This is described in detail in section 4.2. The second system is a supervised classification algorithm, consisting of two parts. First a convolutional neural network operates on lines of text to generate an intermediate representation, which is then combined with the output from the clustering step before being fed into a feed-forward neural network layer, which outputs the probability of each line of text being the start of a new speech. More details on this process are given in section 4.2. Figure 5 shows a high-level overview of the two systems and how they interact.

4.1 Unsupervised

The unsupervised algorithm detects and classifies blocks of text in the PDF file; Figure 6 shows an example of what a desired output could look like. This approach consists of two separate clustering steps, based on work by Klampfl

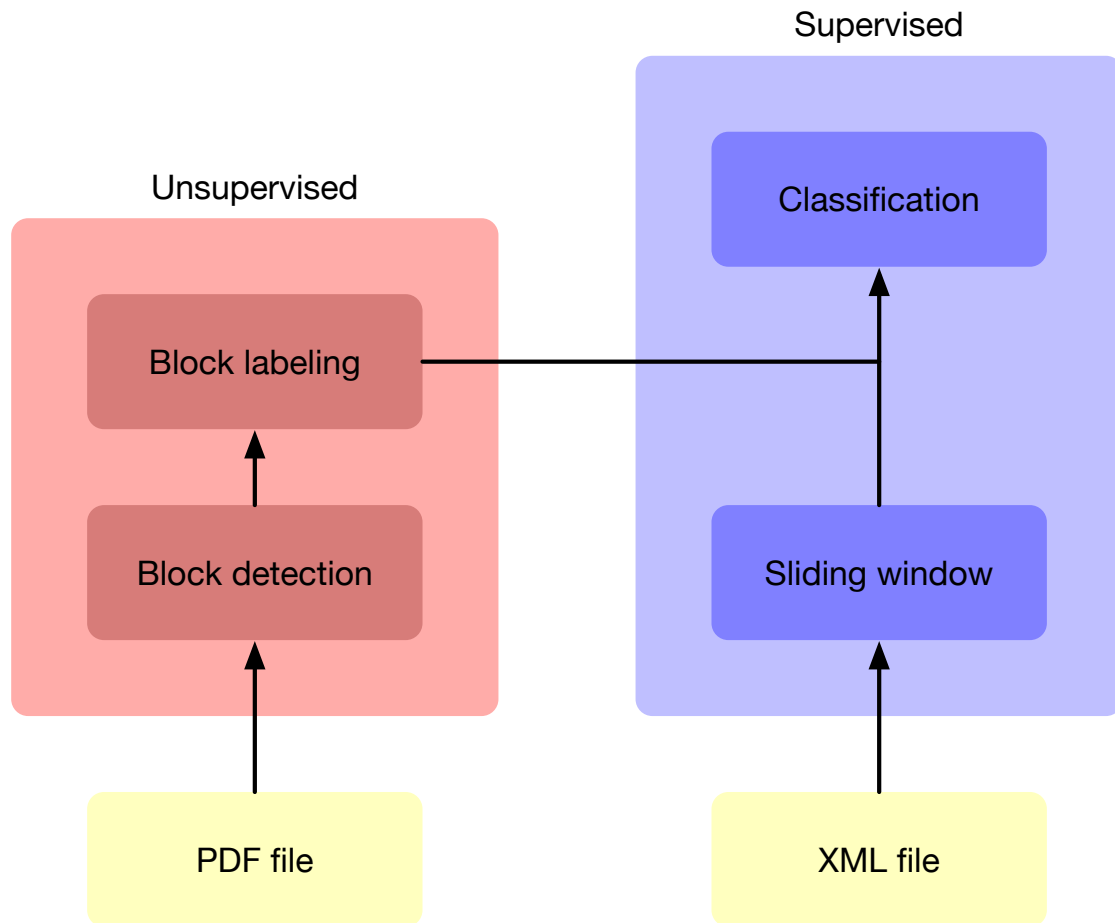
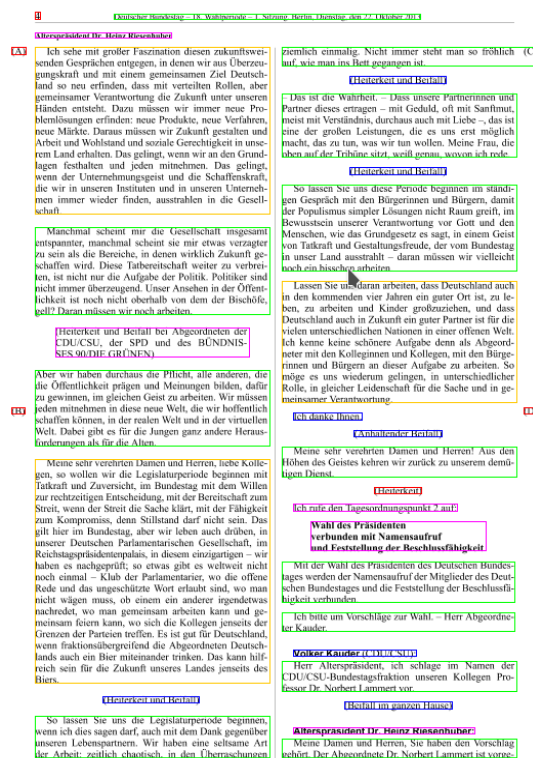


Figure 5 – A high-level overview of the system. The unsupervised block augments the input to the classifier.

213 *et al.* [4]. First, individual characters (the fundamental objects available in
 214 a PDF file) are clustered together into blocks of semantically relevant text.
 215 These could be, for example, paragraphs, section headers or page decoration.
 216 By using the bounding boxes of the blocks, they can be clustered based on
 217 their shape and some additional metadata (e.g. occurrence of font types and
 218 sizes). The rest of this section will go into details on the two clustering steps.

219 4.1.1 Hierarchical Agglomerative Clustering

220 The first step is performed using hierarchical agglomerative clustering (HAC),
 221 an unsupervised bottom-up clustering algorithm that constructs a hierarchical
 222 tree of clusters (in this context referred to as a *dendrogram*). An example
 223 is shown in fig. 7. The algorithm iterates through the individual characters



from the PDF files, successively grouping the two closest clusters (the initial inputs being regarded as clusters of one element) together until only a single cluster remains. This process involves two parameters:

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

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

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

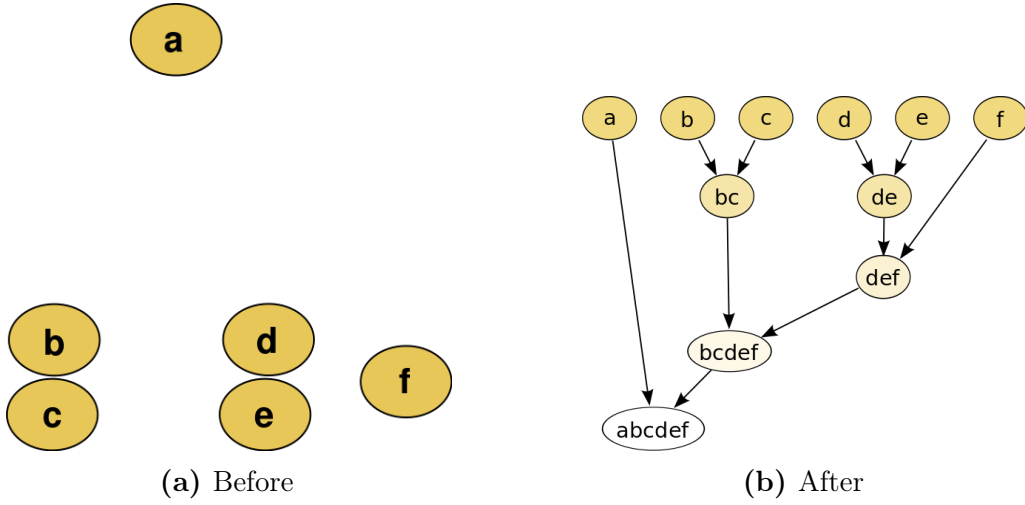


Figure 7 – An example of hierarical agglomerative clustering, where the nodes are clustered by distance.

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

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

More involved linkage criteria exist, such as the Ward method which minimizes the variance within clusters. Although these more complex methods would generally be favored, in this context single-linkage is actually the best option[4]. This is due to its tendency to form long, thin clusters; this mirrors the nature of text, in particular words and sentences (which are just long, thin strings of letters and words respectively). As an additional bonus, it is the most computationally efficient method. While the general time complexity for HAC is in $\mathcal{O}(n^3)$, clever algorithms exist for single-linkage clustering that fall in $\mathcal{O}(n^2)$ [11], making it far more usable on larger datasets.

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

4.1.2 Clustering

The blocks that were found in the previous step are then clustered according to their similarity based on the following metrics:

- Width of the block
- Height of the block
- The most common font occurring in the block
- The size of the most common font occurring in the block

This is implemented two different ways, with their performance to be compared later on. The first method is the familiar K-Means clustering algorithm, which works as follows:

1. Randomly initialize k cluster centroids.
2. Assign each observation to its nearest centroid.
3. Recompute each centroid to be the mean of all of its assigned observations.
4. Repeat steps 2 and 3 until the assignments stop changing.

After the algorithm has converged, the clustering is defined by each observation's nearest cluster centroid; with k clusters, each observation is assigned a k -dimensional one-hot vector.

4.1.3 Dirichlet Process Clustering

K-Means clustering can be generalized to a *mixture of Gaussians* model. Whereas K-Means clustering defines each cluster by a centroid and each assignment by its nearest cluster, a mixture of Gaussians — as the name implies — defines each cluster by a Gaussian distribution and each assignment by a vector indicating the probabilities of belonging to each cluster. This adds a degree of uncertainty to the representation, which is lost in K-Means' hard assignments. While this already improves upon K-Means by increasing the amount of information gained, it still requires a suitable value for the number of distributions, analogous to the K in K-Means clustering. This is, at least in this case, an unintuitive parameter that essentially has to be guessed and evaluated in order to choose a suitable value. Luckily it can be eliminated by using an infinite mixture model. As the name implies, an infinite mixture model is similar to a Gaussian mixture model, except using an infinite amount of distributions, thereby removing the most significant parameter.

The core component of this infinite mixture model is a *Dirichlet process*. This is conceptually similar to the well-known Dirichlet distribution, with one key difference. Whereas the Dirichlet distribution generates discrete probability distributions with a finite amount of possible values, the Dirichlet process generates distributions with an infinite amount of possible values. In a Dirichlet process, the proportions of how many samples are assigned to each cluster are generated using a *stick-breaking process* (alternatively, with a

Pólya urn scheme or a *Chinese restaurant process*)[12]. In the stick-breaking process, the probability w_k of a sample being assigned to the k 'th distribution is given by

$$w_k = \beta_k \cdot \prod_{i=1}^{k-1} (1 - \beta_i), \quad (1)$$

where β_k is a random variable drawn from the distribution $\text{Beta}(1, \alpha)$. Since the Beta distribution is defined on the interval $[0, 1]$, the β_k values can be considered as portions of a unit-length stick. When the first cluster is assigned to, a piece of size β_k is broken off of this unit-length stick. On subsequent assignments, the new sample is either assigned to an existing cluster with a probability proportional to the length of that cluster's portion of the stick, or a new cluster is assigned to. In the latter case, the new cluster gets a β_k -sized portion of the remains of the stick. This way, cluster assignments tend towards a long-tailed distribution. The α parameter in the prior distribution $\text{Beta}(1, \alpha)$ is called the weight- concentration parameter. As shown in fig. 8, the probability mass of the distribution is inversely proportional to the value of α . For lower values of α , big portions of the stick are likely to be broken off, concentrating the cluster assignments into different clusters; conversely, higher values of α lead to smaller portions and more diversity in the cluster assignments.

The tendency to continuously decrease the probability of assigning to a new cluster is a key factor in using Dirichlet process clustering in practice. Since at a certain point the probability of a new cluster being assigned becomes practically zero, it is sufficient to implement this using a finite "large enough" number of clusters, after which the clusters with very low probabilities can be pruned. This is illustrated with a simple example in fig. 9. Note that because of the Bayesian nature of the algorithm, even though the prior favors a long-tail distribution, it has no problem creating equally likely clusters if the data demands it.

The samples generated from each clusters are assumed to come from a Gaussian distribution; the rest of the probabilistic model mostly consists of run-of- the-mill priors, without much semantic meaning and differing between different implementations of the algorithm. In this case an implementation from Scikit- Learn[13] was used, whose probabilistic model and subsequent derivation of the inference algorithm can be found online⁵.

4.2 Supervised

After the data is augmented by the previously described clustering algorithms, it's fed into a convolutional neural network for classification. Since the documents have a dual column layout (meaning each line of text is pretty

⁵<http://scikit-learn.org/stable/modules/dp-derivation.html>

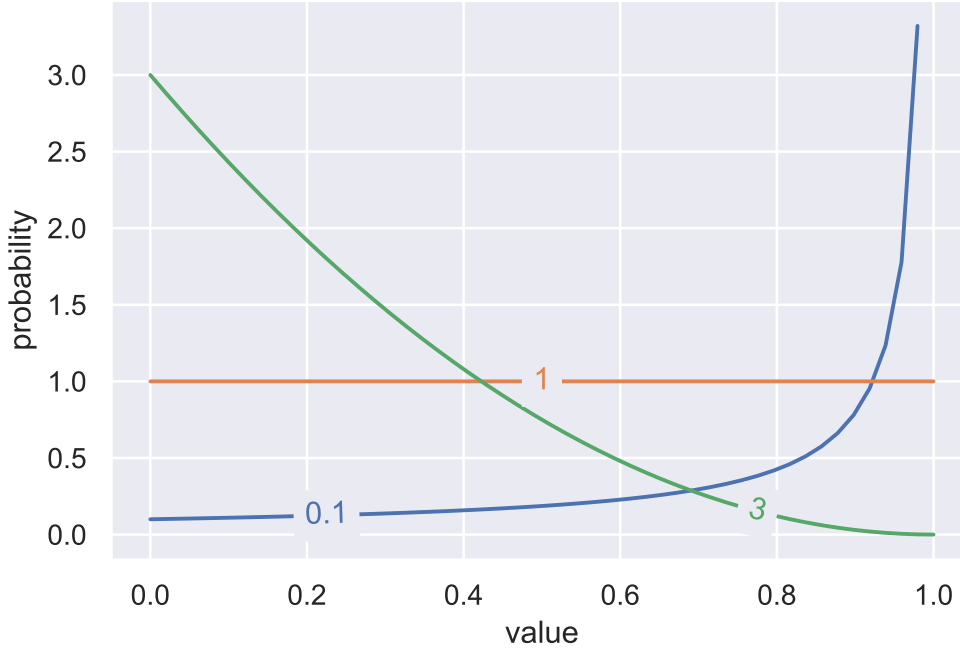


Figure 8 – The Beta($1, \alpha$) distribution for various values of α . As the value of α increases, the probability mass of the distribution shifts towards zero.

short) and classification is based on the lines from the document, a sliding window is used to supply more context. The window consists of a variable amount of lines before and after the main line, which is the one supplying the label (whether or not the line starts a new speech). This sliding window is then used as input to a neural network, consisting of a convolutional part followed by a fully-connected neural network part. After the convolutional part is applied to the sliding window, the output is concatenated to the clustering distribution obtained previously (in the case of the mixture model, this is the actual probability distribution over each cluster; in the K-Means case, it is a one-hot vector). There are two different models, differing in their convolutional layer:

- WordCNN: A shallow word-based architecture[3], where the text is tokenized such that each token is either a word or a single punctuation mark (for example, “Hello-!” would produce the tokens “Hello”, “-” and “!”).
- CharCNN: A character-based architecture[8]. This is a deeper network, that operates on the raw characters without using an embedding layer.

The WordCNN and CharCNN models are described in tables 1a and 1b ;

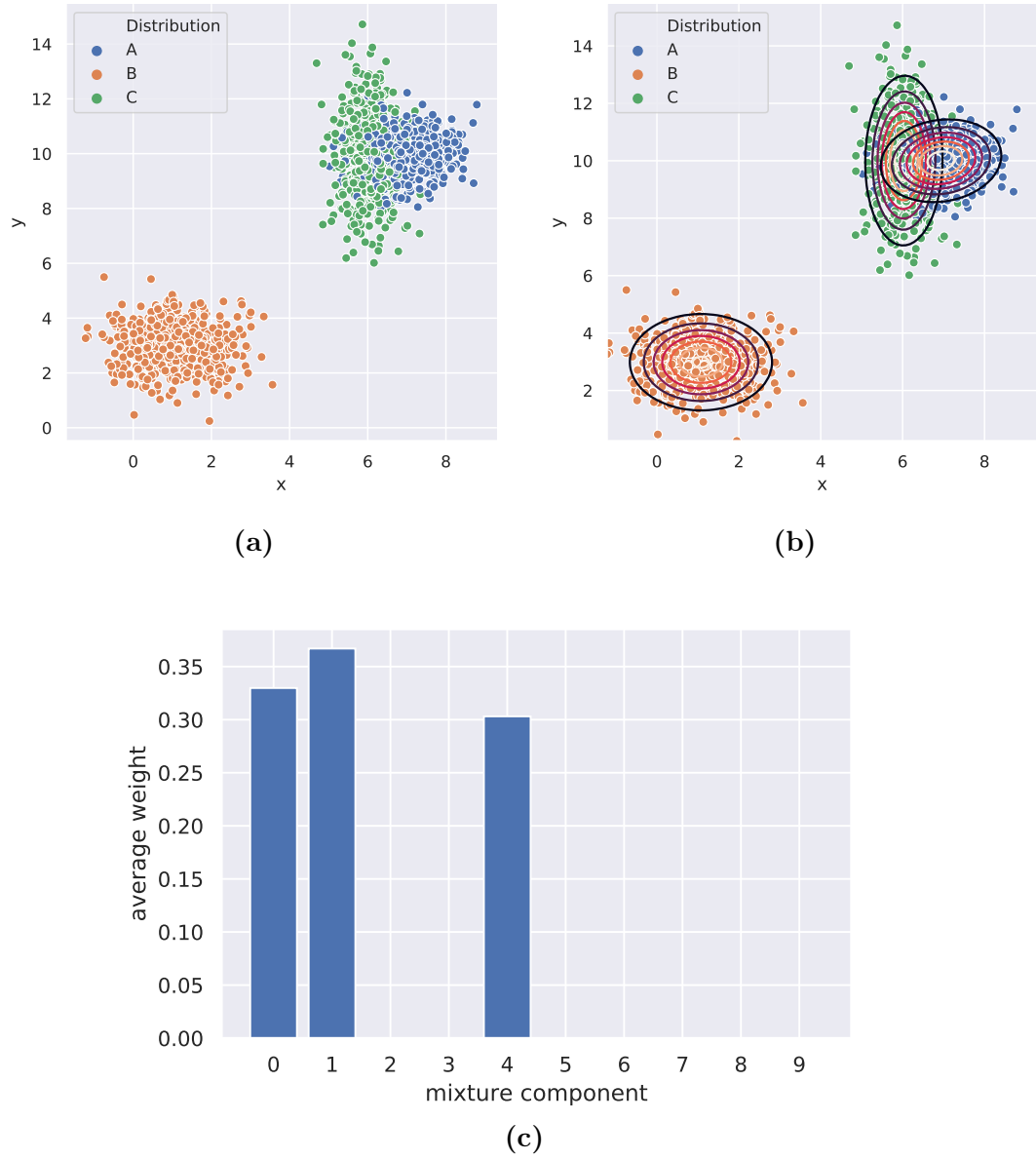


Figure 9 – Dirichlet Process mixture model clustering with a prior of 10 mixture components, performed on 1500 samples drawn from three different bivariate normal distributions (fig. 9a). The distributions found by the clustering algorithm are drawn in fig. 9b. As shown by averaging the component weights of each sample’s assignment (fig. 9c), the Dirichlet Process rightfully assigned to just three mixture components, despite its prior of ten.

Layer	Type	Filters	Size	Pooling
1	Embedding	-	100	-
2	Conv	99	3, 5, 7	1-max

(a) The convolutional part of the WordCNN models consists of an embedding layer to embed the tokens into a higher-dimensional space, followed by a single convolutional layer. The convolutional layer uses 33 filters of each of three different sizes, with a stride of 1. The filters are concatenated for a total of 99 filters, with 1-max pooling applied to each filter such that the final output is a 99-dimensional vector.

Layer	Filters	Size	Pooling
1	256	7	3
2	256	7	3
3	256	3	-
4	256	3	-
5	256	3	-
6	256	3	3

(b) The convolutional part of the CharCNN model is a relatively deeply layered sequence of convolutions. Each convolution is followed by a ReLU activation function. The convolutions have a stride of 1, while the pooling layers are non-overlapping with a stride of 3.

Layer	Output Size	Activation	Dropout
1	1024	ReLU	Yes
2	1024	ReLU	Yes
3	1	Sigmoid	No

(c) Both architectures go through this fully connected neural network with 3 layers. For the layers using dropout, a rate of 0.5 is used.

Table 1 – The layouts of the neural networks. Table (a) corresponds to the convolutional part of the WordCNN model while Table (b) similarly corresponds to the CharCNN model; Table (c) describes the fully connected neural network common to both models.

350 the fully-connected neural network common to both models in described in
351 table 1c. A high-level overview of how the different parts of the system fit
352 together is shown in fig. 5.

353 The network is trained for 100 epochs, stopping early once no significant
354 improvement has been made on a small validation set for 10 epochs in a
355 row. The optimisation process is done using the Adadelata[14] algorithm⁶,
356 with binary cross-entropy as the loss function and the training data delivered
357 in batches of 50 samples. Regularisation is done through a combination of

⁶Although Adam[15] has become very popular, the architecture[3] and survey paper[7] that this work is based on both slightly predate the paper that popularized Adam. Since a full analysis of different optimizers is outside the scope of this thesis and some quick informal tests showed no difference between using Adadelata or Adam, there is no reason to deviate from the survey.

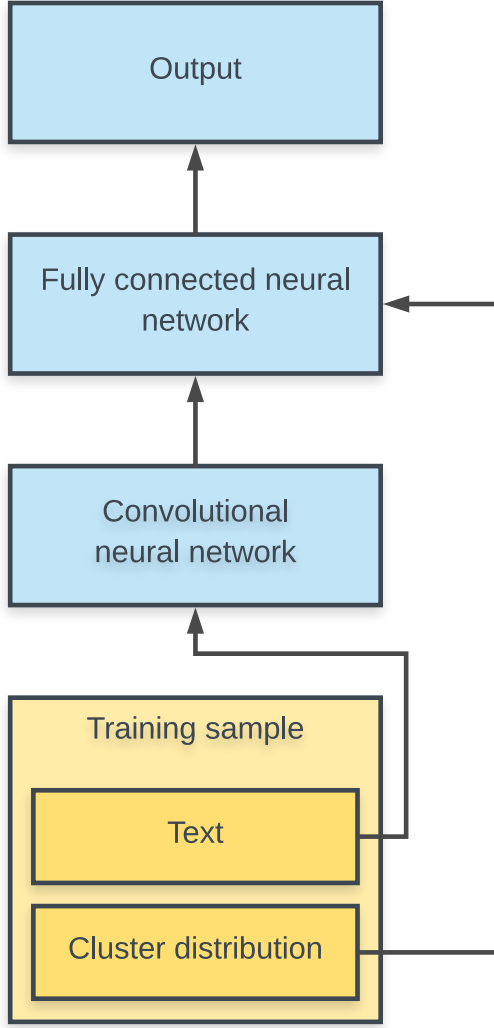


Figure 10 – The layout of the supervised portion of the system. The input data contains text and a cluster type assigned by the unsupervised portion. The text gets put into a convolutional neural network, the output of which is fed together with the cluster type into a fully connected neural network. The sigmoid function is applied to the output of this final neural network to obtain the classification.

dropout[16] and an upper limit on the L2 norm of each weight vector[17]. Further explanations of the optimisation process and the regularisation methods are provided in section 4.2.1 and section 4.2.2 respectively.

4.2.1 Optimisation process

Optimisation is done using the Adadelta update rule[14]. This is easiest to explain by starting with the basic mini-batch stochastic gradient descent algorithm (SGD). At each iteration t , the network parameters θ are updated based on some calculated difference:

$$\theta_{t+1} = \theta_t - \Delta\theta_t. \quad (2)$$

366 The difference between optimization algorithms is how $\Delta\theta$ is calculated. With
 367 SGD, it is simply

$$\Delta\theta_t = \mu \nabla_{\theta_t} \mathcal{L}(\theta_t), \quad (3)$$

where \mathcal{L} is the loss function (in this case, the binary cross entropy between the predicted labels and the true labels), and μ is an arbitrary learning rate between 0 and 1. This learning rate is a tricky parameter to set; too low and learning will take ages, too high and the network will fail to converge because the steps taken are too big. One way to improve this is by using the Adagrad[18] algorithm. While SGD uses a single learning rate for the entire parameter vector, Adagrad (which stands for “adaptive gradient”) adapts the learning rate for each individual parameter; frequently updated parameters get a lower learning rate, while less frequently updated parameters are updated with a higher learning rate. This is done by simply dividing the learning rate with the L2 norm of the sum of all previous gradients:

$$g_t = \nabla_{\theta_t} \mathcal{L}(\theta_t) \quad (4)$$

$$\Delta\theta_t = \frac{\mu}{\sqrt{\sum_{\mathcal{T}=1}^t g_{\mathcal{T}}^2}} g_t. \quad (5)$$

368 This has been found to work very well, in particular in natural language
 369 processing and computer vision where features are often sparse, but it has
 370 two drawbacks:

- 371 1. A suitable value for the global learning rate μ has to be manually
 372 provided.
- 373 2. Since the learning rate is rescaled using a monotonically increasing sum
 374 of previous gradient magnitudes, the learning rate will converge to zero.

375 Adadelta nullifies these drawbacks by eliminating the global learning rate and
 376 restricting the accumulation of gradients to a window of recent updates. First
 377 of all, the sum over all previous gradients is replaced by an exponentially
 378 decaying average of the squared gradients, referred to as $E[g^2]$:

$$E[g^2]_t = \rho E[g^2]_{t-1} + (1 - \rho) g_t^2. \quad (6)$$

Here ρ is a constant representing the rate of decay; this decaying sum serves as a more efficient approximation of an actual window of past gradients, which would require far more memory to store (considering both the huge number of parameters in a neural network and the memory constraints of running on a GPU). Substituting this into the Adagrad algorithm gives us:

$$\Delta\theta_t = \frac{\mu}{\sqrt{E[g^2]_t}} g_t \quad (7)$$

$$= \frac{\mu}{RMS[g]_t} g_t. \quad (8)$$

379 The quantity $\sqrt{E[g^2]_t}$ is called the *root mean square* of g , which occurs often
380 enough in optimisation algorithms that it is often abbreviated as $RMS[g]$. As
381 a final step, the learning rate μ is eliminated by replacing it with a decaying
382 average of the previous gradient updates, similar to eq. (6):

$$\Delta\theta_t = \frac{RMS[\theta]_{t-1}}{RMS[g]_t} g_t. \quad (9)$$

383 4.2.2 Regularisation

384 Dropout[16] is a Regularisation method that is rather crude at first glance;
385 on each batch update, each node in the neural network has a probability
386 p of outputting zero (i.e. the node being “disabled”). As a result, nodes
387 cannot rely on the output of any other node being present, preventing co-
388 adaptation and as a result reducing the probability of overfitting on the
389 training data. Dropout is only active during training; when running the
390 network in evaluation mode, all nodes are active and the outputs are rescaled
391 by a factor of $1 - p$ to account for the now higher activation values. There is
392 another way to view dropout. It is commonly known that neural network
393 (or really any machine learning classifier) performance can be improved by
394 training a large number of them and averaging their outputs. Since every
395 combination of nodes disabled by dropout could be considered a unique
396 neural network, dropout acts as computationally cheap approximation to
397 averaging multiple networks.

398 In addition to dropout, a maximum L2 norm is imposed on each node’s
399 incoming weight vector. If the weight vector’s L2 norm exceeds this limit,
400 it is rescaled so that its new norm is equal to the maximum allowed norm;
401 otherwise the weight vector is left alone:

$$\mathbf{w} = \begin{cases} \mathbf{w} & \text{if } \|\mathbf{w}\|^2 \leq n \\ n \frac{\mathbf{w}}{\|\mathbf{w}\|^2} & \text{otherwise} \end{cases} \quad (10)$$

402 This differs from the more common L2-regularisation — where the combined
403 L2 norm of all weight vectors is added to the loss function — in that the
404 weights are not being continuously pushed towards zero, making it a milder
405 form of regularisation that allows for a bit more complexity in the model.

406 4.2.3 Convolutional Neural Networks

407 Although convolutional neural networks (CNNs) are more closely associated
408 with computer vision, they are also widely used for text classification, offering
409 results competitive with recurrent neural networks[19] at greater speed due
410 to the more parallel nature of CNNs[20].

411 Convolutional neural networks (CNNs) work by taking a number of filters
412 of a specified size and convolving these over the input data. A simplified

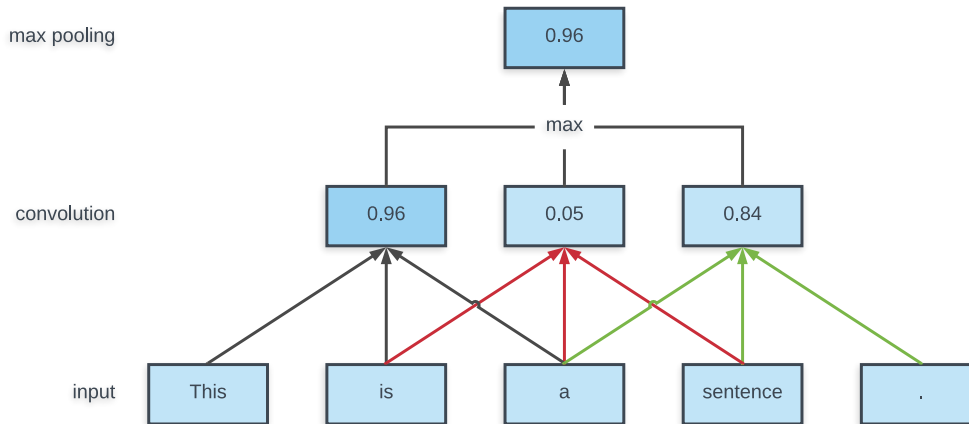


Figure 11 – A simplified convolutional neural network with one filter and max pooling.

example using one filter is shown in fig. 11. In this 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 [21]. 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 superficially 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) or reducing the entire sequence to 1 value (1-max pooling). The reason for this is twofold:

1. It downsamples the data, 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. 11 reacts strongly to the first three words. Without the pooling layer, changing the location of these words in the input would similarly change the location of the high activation in the intermediate representation (*equivariance*). The more aggressively the pooling is applied, the higher the degree of invariance (with full translational invariance being achieved with 1-max pooling).

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

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

5 Experimental Setup

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

1. Number of training epochs until convergence
2. The F1 score or average precision metrics on a test set (see Section 5.2)
3. The number of training samples required to attain a specific F1/average precision score

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

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

5.1 Dataset

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

5.2 Testing performance

All models will be tested on a test set containing 1000 positive samples and 10000 negative samples, all of which are guaranteed not to be in the training

set. Performance on this set is measured by constructing a precision-recall curve, and calculating two values:

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

5.3 CNN performance

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

5.4 Generalisation

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

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

6 Evaluation

To recap the previous sections, we are dealing with two models here:

- “WordCNN”, a convolutional neural network operating on tokenized words and punctuation,
- “CharCNN”, a convolutional neural network operating on tokenized characters,

as well as three variations for each model:

- No clustering.
- The data is augmented with K-Means clustering.
- The data is augmented with clustering through an infinite mixture model.

These variations will be referred to as “Baseline”, “K-Means” and “Infinite Mixture model” respectively. Although K-fold cross validation would seem to be a good way of comparing their relative performance, there are two properties of cross validation that make it less desirable in this case:

- One of the more interesting variables in this context is the size of the training set, but with cross validation it is directly tied to the number of iterations (e.g. with a dataset of 1000 samples and 10 folds, each fold would have a training set of 100 samples).
- The training set and the test set are sampled from the same pool of data. In our case, we want the training set and the test set to be sampled from different electoral periods, because:
 - A change of electoral period means at least a partial change in parliamentary members, which reduces the ability of the model to overfit on the names of members.
 - Although older documents use the same layout as newer ones, the PDF files were generated using different software (electoral period 13 started in 1994 and was the first period to use truly digital documents — older documents are scanned versions of printed paper). Since PDF decompilation is already a flaky process, this results in the PDF decompiler making different mistakes than it does on the newer documents. This is a big problem for rule-based systems, but hopefully the probabilistic nature of neural networks makes them more robust to this issue.

Therefore, a variation of cross validation is used instead. There is a pool of training data consisting of documents from the 18th electoral period of the *Bundestag*, and a test set containing documents from the electoral periods 13 to 17. For 5 iterations, k samples are pulled from the pool of training data. Each of the three models is trained on these samples and then tested on the full test set.

6.1 Parameter Exploration

Before directly comparing the two models, a good baseline value has to be determined for the two parameters in the model that are the most difficult to determine:

1. the size of the sliding window applied to the lines of text in the input data
2. the number of clusters to detect in the blocks of text (i.e. the “k” in k-means)

Although there are more parameters, they are related to the neural network, meaning there is a large amount of prior knowledge available to set the

parameters.

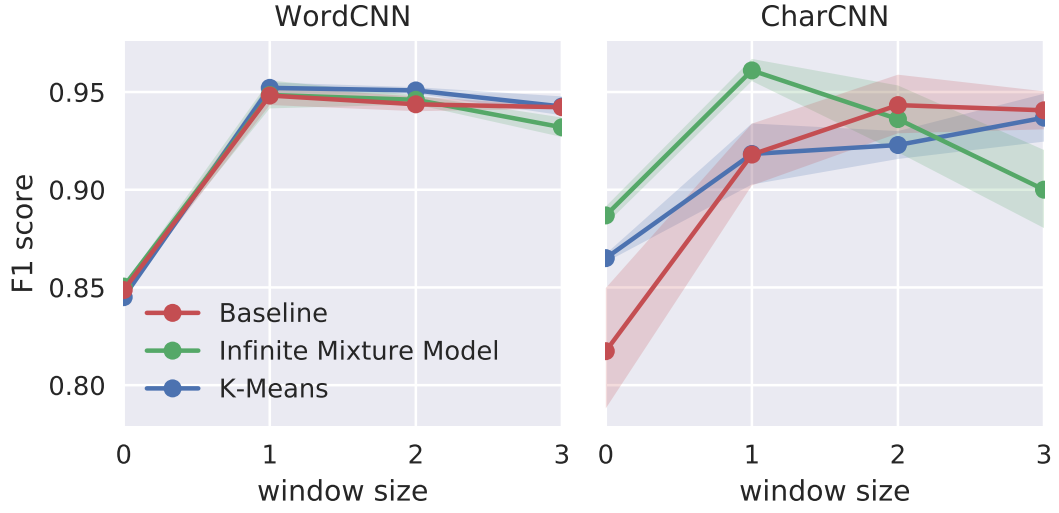
Plot opnieuw genereren met een logischere unit op de x-as, en testen met window sizes die alleen naar voren kijken

The effect of the sliding window size is shown in fig. 12a. In this figure, the window size refers to how many neighbouring lines of text are considered: 0 means just the line of text by itself, 1 means the next line is included, 2 means both the next and previous lines are included, et cetera. For each model, there is a very sharp increase in performance when going from zero to one neighbouring element, followed by a gradual decline as more elements are added.

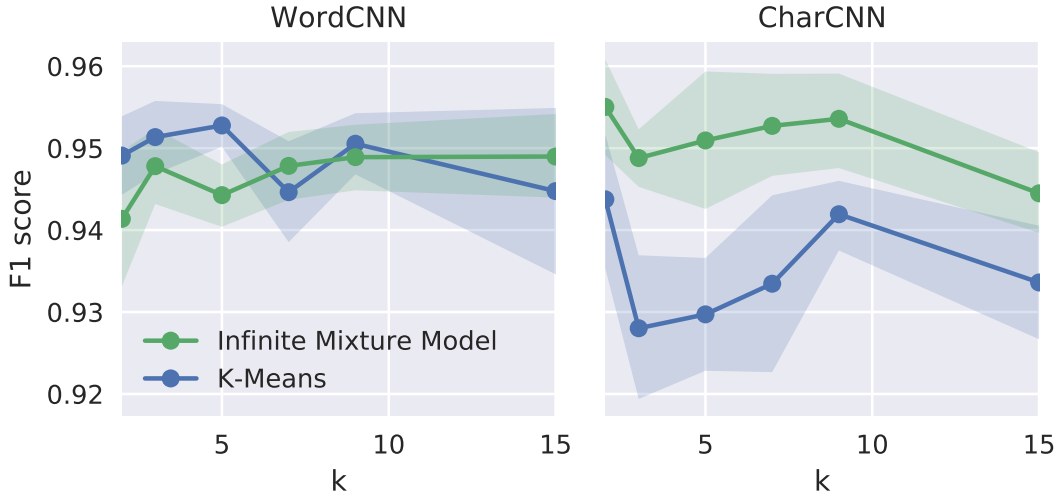
Results for the number of clusters are given in fig. 13. In this case, the window size was fixed to the previously determined optimum of 2. The baseline model is not tested here due to not using clustering, and the model using infinite mixture clustering is not tested due to its fixed amount of clusters. The result varies little, and given the large variances no value can really be said to outperform the rest.

6.2 Training Set Size

Using the ideal parameters obtained in section 6.1, the influence of the training set size is shown in fig. 14. This shows a number of interesting differences. First of all, the Baseline and K-Means models perform similarly for both CNN models at every input size, while the infinite mixture version outperforms them only when using the CharCNN model. Additionally, the WordCNN and CharCNN models respond much differently to variations of the training set's size; WordCNN maxes out its performance at roughly 800 samples, while CharCNN requires around 2000 samples to reach the same performance. However, when augmented using the infinite mixture model, CharCNN peaks at a slightly higher score than WordCNN.



(a) The F1 score with regards to the sliding window size for each model.



(b) The F1 score with regards to the number of clusters for the K-Means model.

Figure 12 – These figures show the performance with regards to the sliding window size and the number of clusters, using a training set of 2000 samples. In fig. 12a, the K-Means model was trained using 9 clusters. Each figure’s F1 score is averaged over 5 trials; the translucent bands around the lines indicates the confidence intervals, meaning that based on the observed F1 scores and assuming normality, the true mean is 95% likely to fall within that interval. The dots on the lines indicate measurements.

Figure 13 – This figure shows the performance with regards to the number of cluster types for each model trained on 1200 training samples with a window size of 5. The vertical axis shows the F1 score, averaged over 10 trials; the horizontal axis shows the number of cluster types considered in the final clustering step. The translucent bands around the lines indicates the confidence intervals, meaning that based on the observed F1 scores and assuming normality, the true mean is 95% likely to fall within that interval. The dots on the lines indicate measurements.

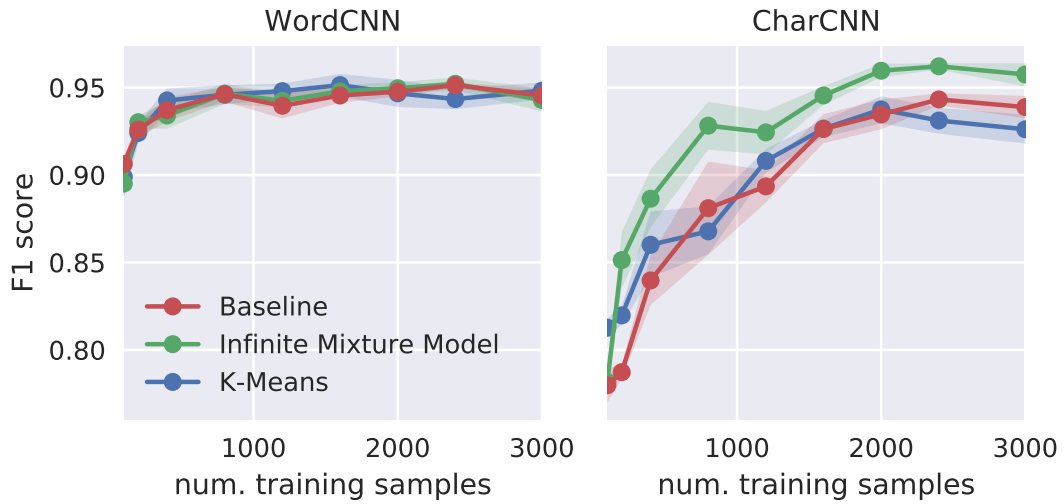


Figure 14 – This figure shows the performance with regards to the size of the training set, using a window size of 2 of 9 clusters for the K-Means models. The F1 scores are averaged over 5 trials; the translucent bands around the lines indicates the confidence intervals, meaning that based on the observed F1 scores and assuming normality, the true mean is 95% likely to fall within that interval. The dots on the lines indicate measurements.

7 Conclusion

Segmenting political proceedings through the use of machine learning is a viable approach, reaching peak F1 scores of 0.96. When training data is very sparse, word-level convolutional neural networks are the best option; augmenting them with information regarding the document layout, whether through K-Means clustering or an infinite mixture model, offers no notable increase in performance. When more training data (roughly 1000 positive samples) is available, character-based convolutional neural networks outperform their word-based brethren only when augmented using infinite mixture model clustering.

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. Klampfl, 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.

- 605 7. Zhang, Y. & Wallace, B. C. A Sensitivity Analysis of (and Prac-
606 titioners' Guide to) Convolutional Neural Networks for Sentence
607 Classification. *CoRR* **abs/1510.03820**. arXiv: 1510.03820. <<http://arxiv.org/abs/1510.03820>> (2015).
608
- 609 8. Zhang, X., Zhao, J. & LeCun, Y. *Character-level convolutional*
610 *networks for text classification* in *Advances in neural information*
611 *processing systems* (2015), 649–657.
- 612 9. He, K., Zhang, X., Ren, S. & Sun, J. Deep Residual Learning for
613 Image Recognition. *CoRR* **abs/1512.03385**. arXiv: 1512.03385.
614 <<http://arxiv.org/abs/1512.03385>> (2015).
- 615 10. Conneau, A., Schwenk, H., Barrault, L. & LeCun, Y. Very Deep
616 Convolutional Networks for Natural Language Processing. *CoRR*
617 **abs/1606.01781**. arXiv: 1606.01781. <<http://arxiv.org/abs/1606.01781>> (2016).
618
- 619 11. Sibson, R. SLINK: an optimally efficient algorithm for the single-
620 link cluster method. *The computer journal* **16**, 30–34 (1973).
- 621 12. Frigyi, B. A., Kapila, A. & Gupta, M. R. Introduction to the
622 Dirichlet distribution and related processes. *Department of Electrical*
623 *Engineering, University of Washington, UWEETR-2010-0006*
624 (2010).
- 625 13. Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python.
626 *Journal of Machine Learning Research* **12**, 2825–2830 (2011).
- 627 14. Zeiler, M. D. ADADELTA: An Adaptive Learning Rate Method.
628 *CoRR* **abs/1212.5701**. arXiv: 1212.5701. <<http://arxiv.org/abs/1212.5701>> (2012).
629
- 630 15. Kingma, D. P. & Ba, J. Adam: A Method for Stochastic Op-
631 timization. *CoRR* **abs/1412.6980**. arXiv: 1412.6980. <<http://arxiv.org/abs/1412.6980>> (2014).
632
- 633 16. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I. & Salakhut-
634 dinov, R. Dropout: A simple way to prevent neural networks from
635 overfitting. *The Journal of Machine Learning Research* **15**, 1929–
636 1958 (2014).
- 637 17. Hinton, G. E., Srivastava, N., Krizhevsky, A., Sutskever, I. &
638 Salakhutdinov, R. Improving neural networks by preventing co-
639 adaptation of feature detectors. *CoRR* **abs/1207.0580**. arXiv:
640 1207.0580. <<http://arxiv.org/abs/1207.0580>> (2012).

- 641 18. Duchi, J., Hazan, E. & Singer, Y. Adaptive subgradient methods
642 for online learning and stochastic optimization. *Journal of Machine*
643 *Learning Research* **12**, 2121–2159 (2011).
- 644 19. Yin, W., Kann, K., Yu, M. & Schütze, H. Comparative Study of
645 CNN and RNN for Natural Language Processing. *CoRR* **abs/1702.01923**.
646 arXiv: 1702.01923. <<http://arxiv.org/abs/1702.01923>>
647 (2017).
- 648 20. Gehring, J., Auli, M., Grangier, D., Yarats, D. & Dauphin, Y. N.
649 Convolutional Sequence to Sequence Learning. *CoRR* **abs/1705.03122**.
650 arXiv: 1705.03122. <<http://arxiv.org/abs/1705.03122>>
651 (2017).
- 652 21. LeCun, Y., Bengio, Y., *et al.* Convolutional networks for images,
653 speech, and time series. *The handbook of brain theory and neural*
654 *networks* **3361**, 1995 (1995).