

Fake news: an algorithmic perspective on fact-checking

Martijn B.J. Schouten
11295562

Bachelor thesis
Credits: 12 EC

Bachelor's degree Information Science

University of Amsterdam
Faculty of Science
Science Park 904
1098 XH Amsterdam

Supervisor
Dr. M. J. Marx

ILPS, IvI
Faculty of Science
University of Amsterdam
Science Park 904
1098 XH Amsterdam

2019-06

Add a
final date

Abstract

Add an
abstract

Contents

1 Introduction

The ability to broadcast information on a large scale has been in the hands of large publishing organizations in the pre-Internet era, but nowadays everyone can share news via social media [?]. This introduces risks on validity and authenticity of news, as social media and digital platforms can speed up the spread of falsehoods without much effort from the author [?].

As a matter of fact, 63% of adults in the United States prefer to read their news on the Internet. Young adults take the lead: 76% of adults between the ages 18 and 49 get their primary news consumption via the web, compared to just 43% for adults of 50 years and older [?]. As time passes by, social media is slowly becoming the primary source of news for more and more people.

The main danger of this development is that human perception is often skewed with regards to objectivity of facts. Naïve realism let consumers of news belief that their perception is right, while other's perceptions are uninformed. Furthermore, confirmation bias results in consumers preferring information that confirms beliefs they already have [?]. This makes consumers vulnerable for the spread of misinformation or fake news.

According to the European Commission, "*disinformation - or fake news - consists of verifiably false or misleading information that is created, presented and disseminated for economic gain or to intentionally deceive the public, and may cause public harm*" [?]. The answer to the problem of fake news as of recently has been to manually fact-check statements on validity, but, as Shu et al. underlines, one of the downsides to this approach is that fake news typically relates to newly emerging, time-critical events. This means the real news may not be fully verified by proper knowledge bases due to a lack of contradicting claims [?]. An automated approach would both help in solving the problem of human subjectivity and the speed at which false information is spread in the current news spreading landscape. Furthermore, such an approach can help human fact-checking by targetting statements that are most likely to be false.

Natural language processing has been in rapid development over the past years. With the releases of OpenAI's GPT-2 model in February of this year and Google's BERT in the autumn of 2018, state-of-the-art pre-trained textual embedding techniques have shown promising results on various classification tasks [?][?]. Although fake news classification has been attempted before [?][?], performance has been rather low. However, these new pre-trained textual embeddings have not yet been used in the fight against disinformation.

This thesis is focussed on the following research question: *what is the performance of combinations of pre-trained embedding techniques with machine learning algorithms when classifying fake news?* This main question will be answered through the results of the following subquestions:

- RQ1** Which way of pooling vectors to a fixed length works best for classifying fake news?
- RQ2** At what maximum sequence length do neural networks hold the highest accuracy when classifying fake news?
- RQ3** How well do neural network classification architectures classify fake news compared to non-neural classification algorithms?

2 Related Work

2.1 Automatic fake news detection

There have been several attempts in the past to create classifiers for automatic detection of lies and fake news. Wang used both linear and neural classifiers to classify statements from the Liar dataset into 6 possible gradations of truthfulness. Furthermore, he added speaker metadata to improve the result of his classifications. Both with and without introducing speaker metadata, the best performing architecture was found to be a convolutional neural network. With an accuracy of 27% without, and 27,4% with metadata on the test set, Wang was able to perform 6,2% and 6,6% better than the majority baseline of 20,8% [?].

From the same dataset, Khurana extracted linguistic features such as n-grams, sentiment, number of capital letters and POS tags to classify the data into 3 labels instead of the original 6 labels. For classification, she used a set of non-neural classifiers. Her best performing classifier, using gradient boosting, obtained an accuracy of 49,03%, which performed around 5% better than the majority baseline of 44,28% [?].

The British factchecking organization Full Fact has developed an architecture that is able to monitor and factcheck statements from the British Parliament and major media outlets in the United Kingdom. It can automatically factcheck the accuracy of statistical claims, for example [?]. For detecting factual claims from texts, the organization uses InferSent, which is a way of transfer learning that has been proved to perform well for the use case of Full Fact [?].

Various tools with regards to fake news detection are also available. Faker Fact is a tool which can classify texts into a set of categories ranging from satire to agenda-driven, the former identifying humorous intent, the latter identifying manipulation [?]. Hoaxy, on the other hand, allows for the visualization of unverified claims through Twitter networks [?].

2.2 Pre-trained textual embeddings

Traditionally, feature representation for text classification is often based on the bag-of-words model, containing linguistic patterns as features, such as unigrams, bigrams or n-grams. However, these approaches completely ignore contextual information or word order in texts, and are unable to capture semantics of words. As a result, classifiers may be unable to correctly identify patterns, affecting the classification accuracy [?].

As an answer to these problems, pre-trained text embeddings have been rising in popularity, both in use and in research. Before classification is possible, text data needs to be transformed into numbers to be able to be interpreted by classifier algorithms. Fundamentally, text embeddings are vector representations of linguistic structures, allowing for usage of text in classifiers. The process of turning words into these embeddings is typically powered by statistics gathered from large unlabeled corpora of text data [?].

In 2017, Vaswani et al. proposed a novel architecture for embedding text data called the Transformer. With the main aim originally being translation from one sentence in some language to another sentence in another language, Transformers are based on an encoder-decoder model. These models take the sequence of input words, convert it to an intermediate representation, after

which the decoder creates an output sequence.

The main strength of the Transformer architecture is its focus on attention to create the intermediate representation. The encoder receives a sequence of inputs, and reads it at once, as opposed to sequentially (either from left to right or from right to left, as humans do). This allows the encoder to learn the context of a word based on all of the surrounding text [?].

Inspired by the Transformer architecture, numerous new text embedding techniques have been developed which give the possibility to classify texts by making use of the intermediate representation of the Transformer model. As shown by the Bidirectional Encoder Representations from Transformer (BERT) model by Devlin et al., these techniques have beaten existing benchmarks in natural language processing, underlining the importance of context in text embeddings [?].

2.3 Pooling

Linear classifiers need data in a two-dimensional shape to be able to perform calculations. In the case of raw text data, sentences in the dataset have variable word lengths, resulting in a different vector length when turning the text into a vector representation. To turn the vector representations into a uniform length, we can either cut off the vectors at a fixed length (*padding*), or we can perform calculations to reduce the length of the vectors (*pooling*).

In computer vision, feature pooling is used to reduce noise in data. The goal of this step is to transform joint feature representations into a new, more usable one that preserves important information while discarding irrelevant details. Pooling techniques such as max pooling and average pooling perform mathematical operations to reduce several numbers into one [?]. In the case of transforming the shape of the data, we can reduce vectors to the smallest vector in the dataset to create a uniform shape.

Scherer et al. compared performance of two pooling operations on a convolutional neural network architecture. The first pooling method extracted maximums and the second one was primarily based on working with averages. They have shown that a max pooling operation is vastly superior for capturing invariances in image-like data [?].

Shen et al. noted that in text classification, only a small number of key words contribute to the final prediction. As a result, simple pooling operations are surprisingly effective for representing documents [?]. Lai et al., Hu et al. and Zhang et al. use a max pooling layer in a (recurrent) convolutional neural network for identifying key features in text classification [?][?][?]. In the case of text classification, max pooling strategies seem to be the most popular.

2.4 Padding

When padding a sequence, a list of sequences is transformed to a specific length. Sequences longer than the desired length will be truncated to fit the requirement, while sequences shorter than the desired length will be padded by a specified value [?]. To fill the sequences, a value of zero is often used. Hu et al. also use zero values for padding their sequences [?].

Apart from controlling the size of the feature dimension, padding has other uses as well. Simard et al. make use of sequence padding for convolutional

id	11044.json
label	pants-fire
statement	The Mexican government forces many bad people into our country.
subjects	foreign-policy,immigration
speaker	donald-trump
speaker_job	President-Elect
state	New York
party	republican
context	an interview with NBC's Katy Tur
mostly_true_count	37
half_true_count	51
barely_true_count	63
false_count	114
pants_on_fire_count	61

Table 1: An example entry in the Liar dataset.

For each statement, the dataset contains an id, a label, a subject, a speaker, the function title of the speaker, the affiliated state and political affiliation, the context of the statement and a vector with a truthfulness history. An example of such a data entry can be seen in table 1. Wang introduced this truthfulness history to boost the prediction scores, as speakers with a track record of lying are expected to have a lower chance of speaking the truth when classifying new statements. However, for our application we are only interested in the statement itself and its corresponding label. Due to cheapness and spreadability, a large amount of fake news is spread over social media [?]. This means author information and metadata will not readily be available in real world circumstances.

The original dataset has been split beforehand into a test, train and validation set. The train set contains 80% of the total amount of statements, while the test and validation set both contain approximately 10% of the statements.

Check
whether
this
number
is still
correct

3.2 Data preprocessing and cleaning

3.2.1 Filtering statements

The original dataset contained statements ranging from 1 sentence to 19. On closer inspection of statements with the high amounts of sentences, it was found that not all statements were processed from source files into dataframes correctly. As a result, records of some different statements were joined together, forming a single string. To combat this, the following regular expression was used to filter those statements out:

```
\\.json\\t(mostly-true|true|half-true|false|barely-true|pants-fire)\\
```

After applying this regular expression, the total amount of sentences in the statements were reduced from a maximum of 19 to a maximum of 11.

6 labels	3 labels	2 labels
true (16.1%)	true (35,3%)	true (55,8%)
mostly-true (19.2%)		
half-true (20.5%)	half-true (20.5%)	
barely-true (16.4%)	false (44,19%)	false (44,19%)
false (19.6%)		
pants-fire (8.19%)		

Table 2: Distribution of labels from the original label distribution when reducing the amount of labels.

3.2.2 Reducing labels

Wang’s main objective was to classify fake news into a fine-grained category of fakeness [?]. For our main research question, we aim to predict whether the statements are fake news or not. This means the statements do not necessarily need to be distinguished into these fine-grained categories. Because of this, the classifiers used to predict fake news in this research will be trained on the original 6 labels, Khurana’s division into three labels [?], and a binary classification. The division from the original 6 labels into the lesser amounts of labels can be seen in table 2. This way, we can better compare performance of pre-trained embeddings to existing research on this dataset.

3.3 Methods

3.4 Applying embedding techniques

As our main research question is focussed on pre-trained word embeddings, the first step in the classification process is to turn the statements of the Liar dataset into vectors. For this purpose, the Flair framework will be used. Flair contains interfaces for turning words into embeddings, built on the PyTorch platform [?][?]. Using Flair, we have access to the following 5 state-of-the-art pre-trained embedding techniques:

- ELMo (Embeddings from Language Models) [?];
- BERT [?];
- Generative Pre-Training (GPT) [?];
- Transformer-XL [?];
- Flair [?].

To apply these embeddings, the Flair framework first requires a sentence object to be created [?]. For dividing the statements into sentences, the `sent_tokenize` function from the `nltk` package is used.

3.4.1 RQ1

3.4.2 RQ2

3.4.3 RQ3

4 Evaluation

4.1 RQ1

Which way of pooling vectors to a fixed length works best for classifying fake news?

The answer to this question will be given by comparing accuracy of a logistic regression with different pooling strategies (max, min, average).

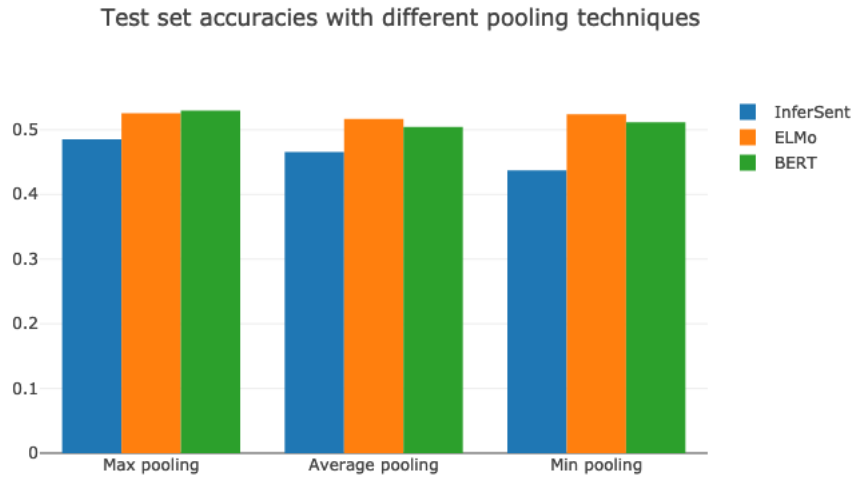


Figure 2: Comparing different pooling strategies.

4.2 RQ2

At what padding sequence length do neural networks hold the highest accuracy when classifying fake news?

The answer to this question will be given by comparing accuracy of a bidirectional LSTM architecture with variable padding lengths.

By comparing all embedding techniques, I plan on getting a peak padding length with the best overall performance shared by all embedding methods. At this moment, the optimal length averages at around a length of 22.

Test set accuracy of padded datasets with variable maximum lengths

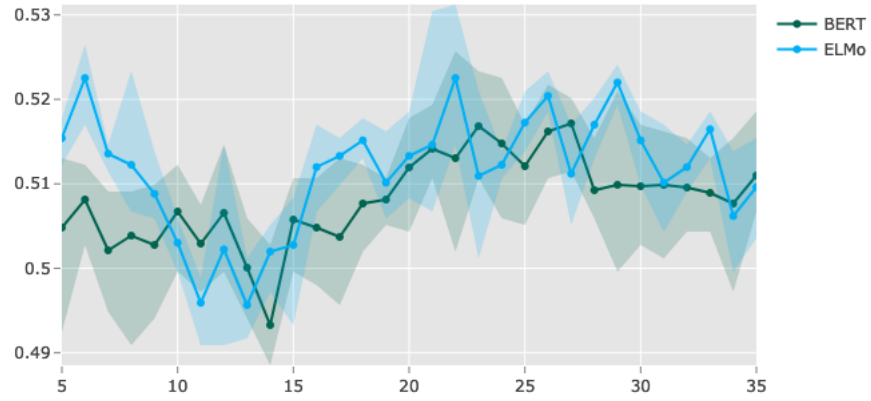


Figure 3: Comparing different maximum padding lengths.

4.3 RQ3

How well do neural network classification architectures classify fake news compared to non-neural classification algorithms?

The answer of this question will be given by comparing two linear classifiers with two neural classifiers.

Test set accuracy of machine learning models for each embedding technique

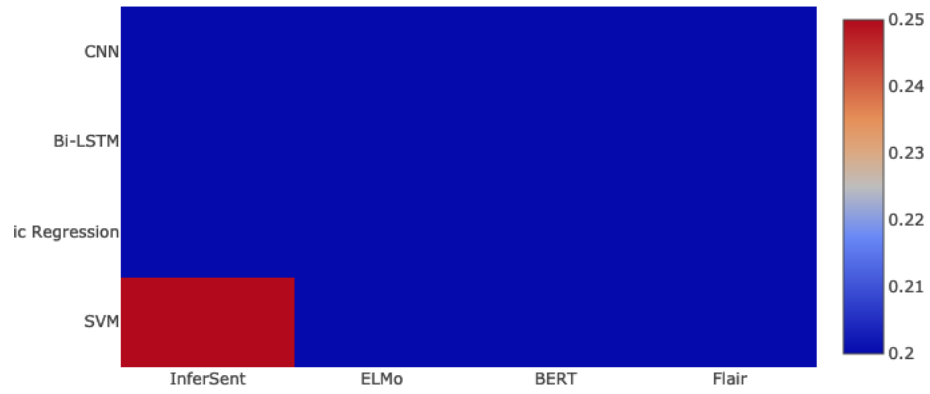


Figure 4: Comparing linear classifiers with neural classifiers.

5 Conclusions

5.1 Acknowledgements