

Word embeddings for predicting political affiliation based on twitter data

Ibrahim Abdelaziz¹, Oliver Berg¹, Angjela Davitkova¹, Md Ashraf Hossain¹, Charu James¹, Shriram Selvakumar¹, Kumar Shridhar¹, and Saurabh Varshneya¹

¹Technische Universität Kaiserslautern

November 21, 2017

1 Introduction

The modern world of social media knows a plethora of means to communicate ones personal opinion and political alignment. With the platform *Twitter*, figures of political interest are expressing their standpoints in small-sized 144-character texts, which contain a comprised message specific to the general public. This yields great potential for automated analysis of party affiliations to classify political persons of interest within the overall political spectrum [3].

Hence do we motivate the structured approach of building a **social media dataset**, utilizing **word embeddings** [7] based modeling approaches to prepare further **qualitative analysis** to obtain dedicated insight and validate possible early intuition. This is to be seen in context of latest **advances in research**.

2 Data assets

Political motives were shown to be consistently predictable with an accuracy better than chance [3]. According to the presented task description of analyzing political affiliation based on Twitter-data, the final comparison occurs on the basis of **tweets** made by political figures on the platform Twitter. Additionally, categorized data taken from *www.wahl.de/politiker* may leveraged this data as a prearrangement of

the initial raw Twitter-data.

For the purposes of this specific research, the main portion of data would be the Twitter data. In order to include all of the relevant politicians from the parties, their corresponding Twitter accounts will be collected (Twitter, wahl.de), where the data is offered as open source data. For each of the Twitter accounts of the politicians, a corresponding party is kept, for which the politician is working. Having these accounts, the same number of tweets will be crawled for each of the parties. 80 percent of the collected data would be used from training the classifier and 20 percent of it would be used for testing. Additionally, as testing data the previously mentioned sources, which incorporate parliament discussion data and party manifesto data, can be used.

3 Model Deployment

In order to train models on text, the data needs to be converted into numerical values, more specifically vectors, under specific similarity metrics. One way this can be achieved is to create a word embedding for each of the words from the tweets using *Word2Vec*. Word embeddings may then apply. It needs to be noted though, that "Word Embedding" is the collective name for a set of language modeling and feature learning techniques in natural language processing (NLP) where words or phrases from

the vocabulary are mapped to vectors of real numbers [1]. Thus, this means that the similar words will have similar vectors and they are going to reside in relatively the same area in the vector space since they both have similar definitions and are both used in similar contexts. While constructing word embeddings, appropriate dimensionality reduction can be applied.

Existing papers [6] then tackle the classification problem using various techniques, such as SVM, SVD or LSM. Mainly these approaches consider solely the political affiliation in America, where the orientation is rather simpler, since there are only two sides and the users are mainly biased towards one side. Other sources tend to focus on comparison of different classifiers when trying to attack this problem, and thus not proposing a complete well-developed approach [2]. Overall, sentiment classification is mostly covered using recurrent- or convolutional neural networks [5].

In connection to the given focus of working on Twitter data, [4] introduces objections to some of the preexisting approaches. With standard classifiers for inferring political orientation having greatly lower accuracy from the accuracy that they report, it is stated that the classifiers cannot be used for classifying users outside the training data. Thus the contradictory arguments hold true.

The proposed solution will leverage neu-

ral networks, where we compare different approaches of both convolutional- and recurrent networks.

4 Analysis of Results

References

- [1] Word embeddings.
- [2] Maneesh Bhandal, Dan Robinson, and Conal Sathi. Text classifiers for political ideologies. 2009.
- [3] Felix Biessmann, Pola Lehmann, Daniel Kirsch, and Sebastian Schelter. Predicting political party affiliation from text. 2017.
- [4] Raviv Cohen and Derek Ruths. Classifying political orientation on twitter: It's not easy!
- [5] Yoon Kim. Convolutional neural networks for sentence classification.
- [6] Arkajyoti Misra and Sanjib Basak. Political bias analysis.
- [7] Maria Pevina1a, Nikolay Arefyev, Chris Biemann, and Alexander Panchenko. Making sense of word embeddings. 2016.