Similarity measure for Public Persons

Andreas Stöckl

andreas.stoeckl@fh-hagenberg.at School of Informatics, Communications and Media University of Applied Sciences Upper Austria Softwarepark 11, 4232 Hagenberg, Austria

Abstract

For the webportal "Who is in the News!" with statistics about the appearence of persons in written news we developed an extension, which measures the relationship of public persons depending on a time parameter, as the relationship may vary over time.

On a training corpus of English and German news articles we built a measure by extracting the person's occurrence in the text via pretrained named entity extraction and then construct time series of counts for each person. Pearson correlation over a sliding window is then used to measure the relation of two persons.

1 Motivation

"Who is in the News!" ¹ is a webportal with statistics and plots about the appearence of persons in written news articles. It counts how often public persons are mentioned in news articles and can be used for research or journalistic purposes. The application is indexing articles published by "Reuters" agency on their website ². With the interactive charts users can analyze different timespans for the mentiones of public people and look for patterns in the data. The portal is bulit with the Python microframework "Dash" ³ which uses the plattform "Plotly" ⁴ for the interactive charts.

Playing around with the charts shows some interresting patterns like the one in the example of Figure 1. This figure suggests that there must be some relationship between this two persons. In this example it is obvious because the persons are both german politicians and candidates for the elections.

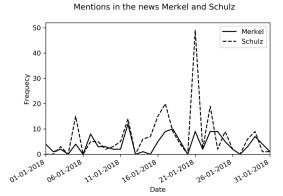


Figure 1: Mentions of Merkel and Schulz in 1/2018

This motivated us to look for suitable measures to caputure how persons are related to each other, which then can be used to exted the webportal with charts showing the person to person relationships. Relationship and distance between persons have been analyzed for decades, for example (Travers and Milgram, 1967) looked at distance in the famous experimental study "the Small World Problem". They inspected the graph of relationships between different persons and set the "distance" to the shortest path between them.

Other approaches used large text corpora for trying to find connections and relatedness by making statistics over the words in the texts. This of course only works for people appearing in the texts and we will discuss this in section 2. All these methods do not cover the changes of relations of the persons over time, that may change over the years. Therefore the measure should have a time parameter, which can be set to the desired time we are investigating.

We have developed a method for such a measure and tested it on a set of news articles for the United States and Germany. In Figure 2 you see how the relation changes in an example of the German chancellor "Angela Merkel" and her opponent on

¹ http://in-the-news.stoeckl.ai/

² http://www.reuters.com/

³ https://dash.plot.ly/

⁴ https://plot.ly/

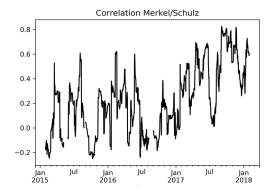


Figure 2: Correlation for Merkel and Schulz

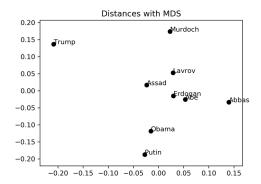


Figure 3: Distances with MDS

the last elections "Martin Schulz". It starts around 0 in 2015 and goes up to about 0.75 in 2017 as we can expect looking at the high correlated time series chart in Figure 1 from the end of 2017.

2 Related work

There are several methods which represent words as vectors of numbers and try to group the vectors of similar words together in vector space. Figure 3 shows a picture which represents such a high dimensional space in 2D via multidimensional scaling (Borg and Groenen, 2005). The implementation was done with Scikit Learn ⁵ (Pedregosa et al., 2011; Géron, 2017; Raschka and Mirjalili, 2017). Word vectors are the building blocks for a lot of applications in areas like search, sentiment analysis and recommendation systems.

The similarity and therefore the distance between words is calculated via the cosine similarity of the associated vectors, which gives a number between -1 and 1. The word2vec tool ⁶ was implemented by (Mikolov et al., 2013b,a,c) and trained

URL	Date	No. Articles
de.reuters.com	2015 to 2018	34058
www.reuters.com	2016 to 2018	36229

Table 1: News articles

over a Google News dataset with about 100 billion words. They use global matrix factorization or local context window methods for the training of the vectors.

A trained dictionary for more than 3 million words and phrases with 300-dim vectors is provided for download. We used the Python library Gensim ⁷ from (Rehurek and Sojka, 2011) for the calculation of the word distances of the multidimensional scaling in Figure 3.

(Pennington et al., 2014) combine the global matrix factorization and local context window methods in the "GloVe" method for word representation ⁸.

(Hasegawa et al., 2004) worked on a corpus of newspaper articles and developed a method for unsupervised relation discovery between named entities of different types by looking at the words between each pair of named etities. By measuring the similarity of this context words they can also discover the type of relatoionship. For example a person entity and an organization entity can have the relationship "is member of". For our application this interesting method can not be used because we need additional time information.

(Zelenko et al., 2003) developed models for supervised learning with kernel methods and support vector machines for relation extraction and tested them on problems of person-affiliation and organization-location relations, but also without time parameter.

3 Dataset and Data Collection

We collected datasets of news articles in English and German language from the news agency Reuters (Table 1). After a data cleaning step, which was deleting meta information like author and editor name from the article, title, body and date were stored in a local database and imported to a Pandas⁹ data frame (McKinney, 2012). The English corpus has a dictionary of length 106.848, the German version has a dictionary of length 163.788.

⁵http://scikit-learn.org/

⁶https://code.google.com/archive/p/word2vec/

⁷https://radimrehurek.com/gensim/

⁸https://nlp.stanford.edu/projects/glove/

⁹ https://pandas.pydata.org/

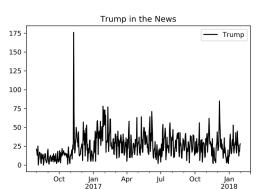


Figure 4: Mentions of Trump

For each article we extracted with the Python library "Spacy" ¹⁰ the named entities labeled as person. "Spacy" was used because of its good performance (Jiang et al., 2016) and it has pre-trained language models for English, German and others. The entity recognition is not perfect, so we have errors in the lists of persons. In a post processing step the terms from a list of common errors are removed. The names of the persons appear in different versions like "Donald Trump" or "Trump". We map all names to the shorter version i.e. "Trump" in this example.

In Figure 4 you can see the time series of the mentions of "Trump" in the news, with a peak at the 8th of November 2016 the day of the election. It is also visible that the general level is changing with the election and is on higher level since then.

Taking a look at the histograms of the most frequent persons in some timespan shows the top 20 persons in the English news articles from 2016 to 2018 (Figure 5). As expected the histogram has a distribution that follows Zipfs law (Adamic and Huberman, 2002; Li, 2002).

From the corpus data a dictionary is built, where for each person the number of mentions of this person in the news per day is recorded. This time series data can be used to build a model that covers time as parameter for the relationship to other persons.

4 Building the Model

Figure 6 shows that the mentions of a person and the correlation with the mentions of another person varies over time. We want to capture this in our relation measure. So we take a time window of n days and look at the time series in the segment

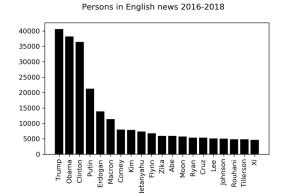


Figure 5: Histogram of mentions in the news

back in time as shown in the example of Figure 1.

For this vectors of n numbers for persons we can use different similarity measures. This choice has of course an impact of the results in applications (Strehl et al., 2000). A first choice could be the cosine similarity as used in the word2vec implementations (Mikolov et al., 2013b). We propose a different calculation for our setup, because we want to capture the high correlation of the series even if they are on different absolute levels of the total number of mentions, as in the example of Figure 7.

We propose to use the Pearson correlation coefficient instead. We can shift the window of calculation over time and therefore get the measure of relatedness as a function of time.

5 Results

Figure 2 shows a chart of the Pearson correlation coefficient computed over a sliding window of 30 days from 2015-01-01 to 2018-02-26 for the persons "Merkel" and "Schulz". The measure clearly covers the change in their relationship during this time period. We propose that 30 days is a good value for the time window, because on one hand it is large enough to have sufficient data for the calculation of the correlation, on the other hand it is sensitive enough to reflect changes over time. But the optimal value depends on the application for which the measure is used.

An example from the US news corpus shows the time series of "Trump" and "Obama" in Figure 6 and a zoom in to the first month of 2018 in Figure 7. It shows that a high correlation can be on different absolute levels. Therefore we used Pearson correlation to calculate the relation of two persons. You can find examples of the similarities of some

¹⁰ https://spacy.io

	Abbas	Abe	Assad	Erdogan	Lavrov	Murdoch	Obama	Putin	Trump
Name									
Abbas	1.00	-0.20	-0.04	0.22	0.21	0.07	0.24	0.20	0.80
Abe	-0.20	1.00	0.27	-0.15	-0.12	0.60	-0.14	0.48	-0.04
Assad	-0.04	0.27	1.00	0.05	-0.03	0.26	0.07	0.24	0.09
Erdogan	0.22	-0.15	0.05	1.00	0.07	-0.02	0.37	-0.25	0.28
Lavrov	0.21	-0.12	-0.03	0.07	1.00	-0.04	0.31	0.17	0.31
Murdoch	0.07	0.60	0.26	-0.02	-0.04	1.00	-0.10	0.80	0.19
Obama	0.24	-0.14	0.07	0.37	0.31	-0.10	1.00	-0.16	0.37
Putin	0.20	0.48	0.24	-0.25	0.17	0.80	-0.16	1.00	0.36
Trump	0.80	-0.04	0.09	0.28	0.31	0.19	0.37	0.36	1.00

Table 2: Similarities of Persons in Dec. 2017

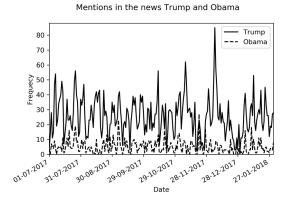


Figure 6: Mentions of Trump and Obama

test persons from December 2017 in Table 2

The time series of the correlations looks quite "noisy" as you can see in Figure 2, because the series of the mentions has a high variance. To reflect the change of the relation of the persons in a more stable way, you can take a higher value for the size of the calculation window of the correlation between the two series. In the example of Figure 8 we used a calculation window of 120 days instead of 30 days.

6 Future Work

It would be interesting to test the ideas with a larger corpus of news articles for example the Google News articles used in the word2vec implementation (Mikolov et al., 2013b).

The method can be used for other named entities such as organizations or cities but we expect not as much variation over time periods as with persons. And similarities between different types of entities would we interesting. So as the relation of a person to a city may chance over time.

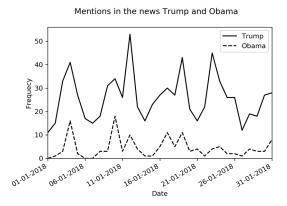


Figure 7: Mentions of Trump and Obama in 1/2018

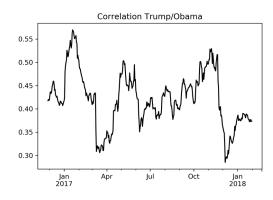


Figure 8: Correlation for Trump and Obama

References

- Lada A Adamic and Bernardo A Huberman. 2002. Zipf's law and the internet. *Glottometrics*, 3(1):143–150.
- Ingwer Borg and Patrick JF Groenen. 2005. *Modern multidimensional scaling: Theory and applications*. Springer Science & Business Media.
- Aurélien Géron. 2017. Hands-on machine learning with Scikit-Learn and TensorFlow: concepts, tools, and techniques to build intelligent systems. " O'Reilly Media, Inc.".
- Takaaki Hasegawa, Satoshi Sekine, and Ralph Grishman. 2004. Discovering relations among named entities from large corpora. In *Proceedings of the 42nd Annual Meeting on Association for Computational Linguistics*, page 415. Association for Computational Linguistics.
- Ridong Jiang, Rafael E Banchs, and Haizhou Li. 2016. Evaluating and combining name entity recognition systems. In *Proceedings of the Sixth Named Entity Workshop*, pages 21–27.
- Wentian Li. 2002. Zipf's law everywhere. *Glottometrics*, 5:14–21.
- Wes McKinney. 2012. Python for data analysis: Data wrangling with Pandas, NumPy, and IPython. "O'Reilly Media, Inc.".
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013a. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013b. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119.
- Tomas Mikolov, Wen-tau Yih, and Geoffrey Zweig. 2013c. Linguistic regularities in continuous space word representations. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 746–751.
- Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. 2011. Scikit-learn: Machine learning in python. *Journal of machine learning research*, 12(Oct):2825–2830.
- Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. Glove: Global vectors for word representation. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543.
- Sebastian Raschka and Vahid Mirjalili. 2017. *Python Machine Learning*. Packt Publishing Ltd.

- R Rehurek and P Sojka. 2011. Gensim–python framework for vector space modelling. *NLP Centre, Faculty of Informatics, Masaryk University, Brno, Czech Republic*, 3(2).
- Alexander Strehl, Joydeep Ghosh, and Raymond Mooney. 2000. Impact of similarity measures on web-page clustering. In *Workshop on artificial intelligence for web search (AAAI 2000)*, volume 58, page 64.
- Jeffrey Travers and Stanley Milgram. 1967. The small world problem. *Phychology Today*, 1(1):61–67.
- Dmitry Zelenko, Chinatsu Aone, and Anthony Richardella. 2003. Kernel methods for relation extraction. *Journal of machine learning research*, 3(Feb):1083–1106.