

Master thesis

Simon Lorent

June 17, 2019

Contents

1	Introduction	4
1.1	What are fake news?	4
1.2	Datasets	5
1.2.1	Fake News Corpus	5
1.2.2	Fake News Net	6
1.2.3	Liar, Liar Pants on Fire	6
1.3	State of the Art	7
2	Data Exploration	8
2.1	Introduction	8
2.2	Dataset statistics	8
2.2.1	Fake News Corpus	8

Master thesis

Fake news detection using machine learning

Simon Lorent

Acknowledgement

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

Master thesis

Fake news detection using machine learning

Simon Lorent

Abstract

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

Chapter 1

Introduction

1.1 What are fake news?

Fake news have quickly become a society problem, being used to propagate false or rumorous informations in order to change behaviors of peoples. Before stating to work on detecting fake news, it is needed to first understand what they are. It have been show that propagation of fake news have had a non negligible influence of 2016 US presidential elections[1]. A few facts on fake news in the United States:

- 62% of US citizen get there news for social medias[2]
- Fake news had more share on facebook than mainstream news[3].

Fake news have also been used in order to influence the referendum in the United Kingdom for the "Brexit".

Shu et al[4] provides multiples features that can be used for fake news detection:

- News Content Features:
 - Linguistic-based
 - Visual-based
- Social Context Features:
 - User-based
 - Post-based

- Network-based

They also provides multiples approaches in order to build the model.

- News Content Models:
 - Knowledge-based
 - Style-based
- Social Context Models:
 - Stance-based
 - Propagation-based

1.2 Datasets

1.2.1 Fake News Corpus

This works uses multiples corpus in order to train and test different models. The main corpus used for training is called Fake News Corpus[5]. This corpus have been automatically crawled using `opensources.co` labels. In other words, domains have been labeled with one or more labels in

- Fake News
- Satire
- Extreme Bias
- Conspiracy Theory
- Junk Science
- Hate News
- Clickbait
- Proceed With Caution
- Political
- Credible

These annotations have been provided by crowdsourcing, which means that they might not be exactly accurate, but are expected to be close to the reality. Because this work focuses on fake news detection against reliable news, only the news labels as fake and credible have been used.

TODO: Expliquer comment le dataset a été nettoyé et mis dans une base de données afin d'augmenter les performances.

1.2.2 Fake News Net

The second dataset used is fake news net[6, 7, 4]. This corpus is made of news from two different sources, PolitiFact and GossipCop. An older version also provides news from BuzzFeed. News are categorized in two classes: fake and non fake. Being quite smaller than fake news corpus, this dataset will be used as a test dataset.

1.2.3 Liar, Liar Pants on Fire

The third and last dataset is **Liar, Liar Pants on Fire** dataset[8], which is a collection of twelve thousand small sentences collected from various sources and hand labeled. They are divided in six classes:

- pants-fire
- false
- barely-true
- half-true
- mostly-true
- true

This set will be used as a second test set. Because in this case there are six classes against two in the other cases, a threshold should be used in order to fix which one will be considered as true or false.

It should be noted that this one differs from the two other datasets as it is composed only of short sentences, and thus it should not be expected to have very good results on this dataset for models trained on Fake News Corpus which is made of full texts.

1.3 State of the Art

Chapter 2

Data Exploration

2.1 Introduction

A good starting point for the analysis is to make some data exploration of the data set. The first thing to be done is statistical analysis such as counting the number of text per class or counting the number of words per sentence. The second step consist of doing Latent Dirichlet Allocation[9] in order to make unsupervised clustering of the text and see if there is some kind of correlation between the clusters to which a text belongs and its labels.

2.2 Dataset statistics

2.2.1 Fake News Corpus

Because **Fake News Corpus** is the main dataset, the data exploration will start with this dataset. And the first thing is to count the number of items per class.

Because the dataset have been cleaned, numbers provided by the dataset creators and number computed after cleaning will be provided. We found the values given at **Table 2.1**. It shows that the number of fake news is smaller by a small factors with respect to the number of reliable news, but given the total number of items it should not cause any problems. But it will still be taken into account later on.

To have a better view of the distribution of categories, an histogram is provided at **Figure 2.1**.

Type	Provided	Computed
Fake News	928,083	
Satire	146,080	
Extreme Bias	1,300,444	
Conspiracy Theory	905,981	
Junk Science	144,939	
Hate News	117,374	
Clickbait	292,201	
Proceed With Caution	319,830	
Political	2,435,471	
Credible	1,920,139	

Table 2.1: Number of texts per categories

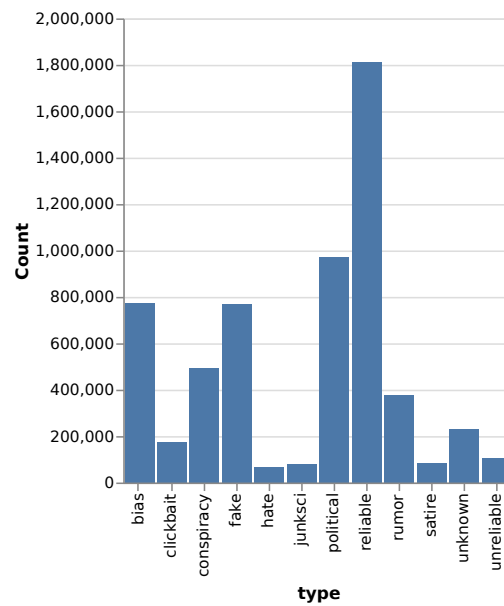


Figure 2.1: Histogram of text distribution along their categories on the computed numbers.

Bibliography

- [1] Hunt Allcott and Matthew Gentzkow. Social media and fake news in the 2016 election. In *Journal of Economic Perspective*, volume 31, 2017.
- [2] Jeffrey Gottfried and Elisa Shearer. *News Use Across Social Medial Platforms 2016*. Pew Research Center, 2016.
- [3] Craig Silverman and Lawrence Alexander. How teens in the balkans are duping trump supporters with fake news. *Buzzfeed News*, 3, 2016.
- [4] Kai Shu, Amy Sliva, Suhang Wang, Jiliang Tang, and Huan Liu. Fake news detection on social media: A data mining perspective. *ACM SIGKDD Explorations Newsletter*, 19(1):22–36, 2017.
- [5] Maciej Szpakowski. Fake news corpus. <https://github.com/several27/FakeNewsCorpus>. Accessed: 2018-10.
- [6] Kai Shu, Deepak Mahudeswaran, Suhang Wang, Dongwon Lee, and Huan Liu. Fakenewsnet: A data repository with news content, social context and dynamic information for studying fake news on social media. *arXiv preprint arXiv:1809.01286*, 2018.
- [7] Kai Shu, Suhang Wang, and Huan Liu. Exploiting tri-relationship for fake news detection. *arXiv preprint arXiv:1712.07709*, 2017.
- [8] William Yang Wang. ”liar, liar pants on fire”: A new benchmark dataset for fake news detection.
- [9] David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022, 2003.