

Stance Detection for Fake News Analysis

First A. Author^{1,2,*}, Second B. Author², and Third C. Author³

¹*Affiliation*

²*Affiliation*

³*Affiliation*

*Contact: First.Author@abc.com, phone +31-50-363 4074

Abstract— In this paper, the implementation of stance detection for the fake news challenge has been discussed and analysed. The methodology has been explained and the experimental results have been discussed.

I. INTRODUCTION

The idea of counterfeit messages in social attitudes is often referred to as a "taste" and is defined as "a collusion to trick you into clicks", Tavernzi (2016). Some news articles contain addresses that interest the reader. However, the author emphasizes only a certain part of the article in the title. If the article itself does not focus or give much truth to what is written in the title, the news can be misleading. The task was performed in Python to find false messages from the registry. Wrong messages were detected using algorithms to find the mood of the data. For this implementation, the mood was determined by analysing the mood of the data. The introduction of the neural network has always been part of a false news mission, and this has been achieved.

The aim of this project is to use natural language processing techniques to automate site detection, as it is not suitable for people to review all information generated by the media.

The actual issue of site detection is to identify the position of the text content in relation to a particular question: the mood of the given text content (positive), negative (negative), or not (neutral) towards that topic. With the concept of interest, we develop a two-stage solution. In the first stage we classify the self - whether the data is neutral or subjective in relation to the subject. In the second stage, we classify feelings of personal data

(ignoring neutral data) - regardless of whether this personal data has a preference or position on the subject. We offer each phase a long-term memory (LSTM) based on a deep neural network that attracts attention to each stage.

There are many ways you can try to identify fake or biased messages on the Internet. However, we believe that our application-based detection mode provides maximum flexibility and reliability without falling into weeds that are true or false for individual needs. Instead, we seek a more general approach by categorizing articles from unknown sources that usually conform or disagree with known power sources (high and low).

In addition, our implementation is particularly convincing because we can accept user input as a reference to an OR article, since any statement needs to verify the truth (Obama is not a US citizen). Therefore, our program acts as a search engine to find facts and return links to related articles, along with the position of this article (agree / agree / neutral) on this topic! Our program offers great opportunities for search and exploration for users, as well as easy verification of claims.

We wanted to create an easy-to-use system to detect the validity of a claim or component of the user based on the concept of site detection. It is difficult to identify erroneous messages. Many "facts" are too complex and difficult to verify or exist in "continuity of truth" or complex sentences with overlapping facts and fantasies. The best way to address this problem is not to check the facts, but to compare how respected sources relate to the claim.

II. LITERATURE REVIEW

Site recognition is the method used to determine the quality of a news article, taking into account what other organizations are writing on the same address. The text of the text must be true, not agreeable, discussed, or not related to the title "Challenge of Fake News" (2016). Site Definition is the method used to determine the quality of the news source. A record is provided by [FakeChallenge.org] (<http://fakenewschallenge.org>), which consists of header and text. This text can come from another article. If texts are allowed from different articles, this system can take into account that other organizations are talking about the same title. The result of the system is the position of the object associated with the name. As described in Fake News, the system will support the following types of placements:

- Approves
- disagree
- Discuss
- irrelevant

Using this system, you can collect article statistics to collect titles. Based on these statistics, the user can come to his own conclusion about whether the new institution has important news sources. To achieve these situations, this system trains the data provided by the mission of the counterfeit message. These data provide a position with the head and body so that the system can detect any combinations of the word leading to the position. For testing, data is provided without articles. To expand the basic level, this project will cover names, stop ending words, and smooth.

After the investigation, I quickly discovered that there are many different categories of misinformation. There are articles that are obviously wrong, articles that present a real event, but then there is some misinterpretation of articles that are pseudo-scientific articles masquerading as news articles that are in fact satirical points, articles consisting mainly of statements and quotations from other people. I came together and found that some people were trying to categorize websites into groups like "spelling," "fake," "misleading," etc.

Information health has become a long-term problem affecting both business and society for print media and digital media. In social networks, the magnitude and consequences of information dissemination are so rapid that distorted, inaccurate or false information has

enormous potential to reach the real world for a few minutes to affect millions of users. Recently, there have been some public concerns about the issue and some methods to alleviate the problem.

Identifying the mood with user content was a long-term problem [9]. However, if the capture position is the user's mood (view) is uncommon, but with respect to a particular subject, it has attracted the interest of researchers only recently. The main work of Mohamed et al. [8] and later on Data 2016 [7], led the authors to the forefront of intensive research in this field.

Data 2016 studied various proposed models, including traditional approaches to computer learning, genetic algorithms and approaches to deeper realization of how the neural network (CNN), the repetitive neural network (RNN) and the long-term memory competition (LSTM). Miter [14] offers the best solution to conduct a thorough investigation into the remainder of the encodings Configuring 256-word dimensions using the word2vec algorithm [6] Skip count, then the second level is the 128 LSTM module. Among other things, pkudlab [12] and DeepStance [11] use CNN's deep models. Augenstein et al. [1] Use a bi-directional attention model.

In some works, a two-stage approach was used. ECNU [15] is determined in the first stage, whether these data relate specifically to the object of the objective, and in the second - determines the orientation (valid / counter). Ltl.uni- because [13] a two-level method of aggregation using SVMs is also used. Among other things TakeLab [2], automated learning mixed with genetic algorithms. Other approaches such as [4] CU-General Women's Union and IUCL RF [5] use the traditional learning machine. Recently, a common problem has been proposed [10].

The mean average results obtained by the participants ranged from 46.19 at the low end to 67.82 at the high end. Recent work has been reported by Du et al. [3], the first of its kind that has deep roots in architecture and uses modelling models. Surpassed the intensive training approach and achieved an F value of 68.79% compared to the F-rating level of the deep training level of 67.82%. We also note that Data 2016 tasks have been classified

as an average F-point of support and return only, ignoring each category (neutral). However, we maintain the measurement accuracy of all three categories (in addition to the F score we measure according to traditional literature) and demonstrate that our model outperforms the deep learning system, not only two classes on average F-score, but also for accuracy Full dimensional measurement.

Despite the fact that in recent months the media in terms of wrong news has been reported, the latest false news problem is rooted in a long history of information operations and disinformation campaigns.

In the detailed paper on information operations, Facebook defines "[...] the actions taken by organized entities ... distorting internal or external political sentiment, which often adds to strategic and / or geopolitical outcomes ..." And "false news" as a useful tool in a series of information operations, the document states that "false news" is "[...] news articles that are supposed to be realistic but contain intentional distortion of the facts in order to stir emotions and attract viewers, to deceive. "

Instead the article declares victory wire, and describes the false news in a simpler way: "[...] through the news to turn social page share media views created that advertising dollars and possibly traction political."

This is clearly a complex problem of cyberspace solution, especially in a world where technology and social media can help make these stories accessible to a wider audience. This has led many researchers in the field of science and industry to develop the FNC. The goal of describing themselves from the FNC is to "[...] solve the problem of false news, to organize a competition to help identify deception messages and deliberate disinformation in the news develop tools to strengthen human inspectors."

The first repetition of the task (FNC 1), which lasted from 1 December 2016 to 2 June 2017 has a position focused exclusively on acquisition, which is the first step in determining the fake news.

While the actual characterization of the truth is a difficult task, full of political and technical problems, finding sites is the first step toward a more reliable

solution. Dean Pomerlo (Dean Pomerlo), one of the organizers of the call, explained the media shift in an interview that "[...] targeting [sites] - to determine who has the best argument, not just the most popular quotes or read extensively how the search engine works. "

Part of the FNC disclosure position can also be defined as a symbol of the relationship that the article text in its title / claim - in particular, whether the body does not agree, or to discuss the title / claim, or if the body is not completely connected. Thus, four possible systems to determine the position of the house must correspond to "Compatibility", "Not OK", "Discussion" and "Offline."

Most studies on natural language processing revolve around research, especially corporate search. This indicates that users may request records in the form of a question they may represent to another person. The device interprets important elements in the human sentence, for example, such as those that may correspond to certain functions of the record and return an answer.

NLP can be used to interpret and analyse free text. There is a large amount of information stored in free text files, for example, patient records. Before NLP models, which relied on in-depth training, this information was not available for computer analysis and cannot be systematically analysed. But NLP allows analysts to search the vast text of free text to find the necessary information in the files.

Mood analysis is another major option to use NLP. Using emotion analysis, data scientists can evaluate social networking comments, such as how their brand works, or comment on customer service group feedback to identify areas where the company should work better.

Google and other search engines rely on their automated translation techniques in NLP's detailed training modules. This allows algorithms to read text on a web page, interpret its meaning and translate it into another language.

III. METHODOLOGY

Stance detection can be formulated in different ways. In the context of this task, we define stance detection to mean automatically determining from text whether the author is in favour of the given target, against the given target, or whether neither inference is likely. Consider the target--data pair:

The data provided is (headline, body, stance) instances, where stance is one of {unrelated, discuss, agree, disagree}. The dataset is provided as two CSVs:

train_bodies.csv

This file contains the body text of articles (the articleBody column) with corresponding IDs (Body ID)

train_stances.csv

This file contains the labeled stances (the Stance column) for pairs of article headlines (Headline) and article bodies (Body ID, referring to entries in train_bodies.csv).

Distribution of the data

The distribution of Stance classes in train_stances.csv is as follows:

rows	unrelated	discuss	agree	disagree
49972	0.73131		0.17828	
	0.0736012		0.0168094	

People can infer from the data that the speaker is likely to be against the target. The task is to test automated systems to see if they can see the registry location. To successfully locate a site, automatic systems often need to identify corresponding pieces of information that may not be present in the focus text. For example, if one supports the rights of the foetus effectively, it is likely that he or she is against the right to abortion. We provide scope limits that pertain to each purpose that systems can collect information to help locate.

Automatic location is often used to retrieve information, compile text, and display text. In fact, it can be said that finding a position can add additional information to the analysis of emotions, because we often appreciate the author's assessment of certain goals and proposals, not just whether the speaker is angry or happy.

The current issue of site detection is to identify the position of the text content on the specific topic: Does this text content mean the value of FAVOR (positive) or locked (negative) or not (neutral) in the direction of this topic. With the concept of interest, we develop a two-

stage solution. In the first stage we classify the self - whether the data is neutral or subjective in relation to the subject. In the second stage, we classify feelings of personal data (ignoring neutral data) - whether the personal data given to them are FAVOR or OBJECT in relation to the subject. We offer each phase a long-term memory (LSTM) based on a deep neural network that attracts attention to each stage. We have achieved the best macro result of 68.84% and the best accuracy of 60.2%, bypassing solutions based on in-depth training. Our structure, T-PAN, is the first in the thematic literature on site detection that uses in-depth training in a two-stage structure.

We then used a set of keywords to gather up to several thousand items from the event log to access the automated learning model. Here we went to collect more than a few articles, because automated learning will determine relevance in the future.

After calculating many APIs for processing newspapers and natural languages, we found that keyword searches are best with relevant articles. The task was to provide a natural language processing algorithm that would extract the most relevant keywords to be found and retrieve just the right number of keywords. I have summed up many algorithms and prepared more than 50 keywords, making them too much to search for. In addition, many algorithms were draining resources, and sometimes it took a minute to analyse that text.

We have created and implemented an automated learning model in Tensorflow, based on several research in the field of positioning [1] [2] [3]. Our model uses a set of word words, Word 2 Vec, TF, TF-IDF and Stop Words in Scikit-learn to express our input. This is done over a hidden level with ReLU activation, a fully connected level, and a softmax activation function to create one of the four outputs. We compare any text with any request. Thus, our ML derives whether our text is "binding" or "unrelated" to the request. If linked, it will be displayed if the company "agrees", "does not agree" or "is neutral" with respect to our claim. Our model achieved a 82% accuracy in test data (to detect a clean rack ... not necessarily "fake news").

Some teams are trying to teach machine models to teach groups of "fake" articles and collections of "real" articles. This method is terrible because fake messages

can appear in articles written correctly and vice versa! The method is not limited to content and we are interested in finding real content.

Some teams try to validate each fact in the article. This is interesting and could eventually become part of a future system to detect counterfeit messages, but today this method is impossible. The fact of the facts is in a continuum and depends to a large extent on the nuances of individual words and their connotations. It is difficult to divide the nuances of the human tongue into real / false divisions.

In order for our application to work, we had to compare the new positions with our constantly updated reputation database. We wrote a Python script to track all sources that were found along with weight. At first, we built a reputation based on research at the state level, and every time we ran our algorithm, we added new sources to our database. To do this, we calculated the degree of reputation for each new article by comparing its status to the input request with source sites for the known reputation and the result rate. In the future, we hope to provide more accurate ways to science-based knowledge to improve our database. Because we are a smaller project, we also hope to find an improved way to track the database using csv, since a copy of the database is located outside the same application start-up.

The structure of this model was chosen because of the ease of execution and the quick calculation because we can rely on the twig instead of the repetition. Judging from the relative strength of this model, turns seem to

capture a wide range of subjects. However, the form is limited to receiving only one note of text. One possible extension of this model is to include some of the recurring interest mechanisms of the package, allowing for specific aspects of the title / body to be modelled after obtaining a general summary of CNN.

IV. CONCLUSIONS

In this paper the implementation of stance detection for the fake news analysis has been made, and the methodology for the detection has been explained with experimental results.

ACKNOWLEDGMENT

The heading of the Acknowledgment section and the References section must not be numbered.

REFERENCES

- [1] Amarasingam, A. (Ed.). (2011). *The Stewart/Colbert effect: Essays on the real impacts of fake news*. McFarland.
- [2] Williams, B. A., & DelliCarpini, M. X. (2011). Real ethical concerns and fake news: The Daily Show and the challenge of the new media environment. *The Stewart/Colbert effect: Essays on the real impacts of fake news*, 181-192.
- [3] Conroy, N. J., Rubin, V. L., & Chen, Y. (2015). Automatic deception detection: Methods for finding fake news. *Proceedings of the Association for Information Science and Technology*, 52(1), 1-4.
- [4] Riedel, B., Augenstein, I., Spithourakis, G. P., & Riedel, S. (2017). A simple but tough-to-beat baseline for the Fake News Challenge stance detection task. *arXiv preprint arXiv:1707.03264*.