

Stance Detection & Fake News NLG Proposal

Anonymous ACL submission

1 Introduction

The term *Fake News* has gained popularity following the 2016 United States presidential election and the vote for the United Kingdom to exit the European Union (Rose, 2017)(Kucharski, 2016). Fake News refers to articles that meet poor journalistic standards, and contain incorrect or misleading information. It's suggested that these articles, and their tendency to be shared on social media had discernible effect on the events of the USA election, and Brexit (Allcott and Gentzkow, 2017).

Determining whether or not a news article is fake is difficult. A Stanford study shows students from middle school through college have difficulty distinguishing real news articles from advertisements (Wineburg et al., 2016). As a result, attempts have been made to automate the detection of fake or misleading news articles (Conroy et al., 2015)

The purpose of this research is to apply natural language processing and machine learning techniques to analyzing the validity of news articles. In particular, we will begin by following the outline presented by Fake News Challenge (Challenge).

The *stance* of a text is the attitude it expresses towards a particular target (Augenstein et al., 2016). The first step of the Fake News challenge is to categorize the stance of the body against the stance of the heading of the article. The Fake News Challenge organization provides an implementation, which we will use as our baseline.

2 Previous Work

2.1 Data Requirements

Rubin et. al. showed data used to investigate rumors and deception need to have the following characteristics. There must be both truthful and deceptive news within the data set, the format must be accessible, the data must be verifiable, and there must exist data points of comparable lengths and writing styles(Rubin et al., 2015).

2.2 Stance Detection

In Emergent: a Novel Data Set for Stance Classification (Ferreira and Vlachos, 2016), stance detection was used to classify claims in relation to news articles in the Emergent data set. Emergent is a data-set derived from a digital journalism project for rumour debunking. Consisting of 300 claims and 2,595 associated news articles, the Emergent project contains labelled data that can be used in a variety of NLP tasks. (Ferreira and Vlachos, 2016). Headlines could be classified as either for, against or observing a claim, where observing a claim merely mentions the claim without giving a stance.

A headline of a news article and an associated claim were considered. Certain features were extracted, considering first the headline alone and then the headline and claim together. Headline features were extracted using the bag-of-words representation, as well as whether the headline ended in a question mark. The idea of the bag-of-words representation is to quantize each extracted key point into one of visual words, and then represent each image by a histogram of the visual words. For this purpose, a clustering algorithm (e.g., K-means), is generally used for generating the visual words /citezhang2010understanding. Features for the distance from the root word to

any refuting words (e.g deny) and any reporting words (e.g claim, presumably) were also added.

For considering the headline and the claim together, their approach involved constructing a graph, such that each word in the headline and the claim was vertex. For every word pairing between the claim and the headline an edge was created and assigned a score. If the stems of the words are identical, it is given the max score. If the words are paraphrases, which is determined using the Paraphrase Database (PPDB) (Ferreira and Vlachos, 2016), they are given the maximum paraphrase score. If neither they are given the minimum score. The Kuhn-Munkres algorithm was run on the graph, which found the maximum scoring word alignment (Ferreira and Vlachos, 2016). Features were added for negating words and subject-to-verb objects. Vectors were used for comparing the claims to headlines, looking at the cosine similarity of the subject-to-verb objects.

Using the headline and headline-claim methods, two overlap thresholds were defined, minimum for and maximum against. If the overlap is higher than the minimum for threshold it is labelled for, if it is less than the maximum against it is labelled against. If it hits neither threshold it is labelled observing. An accuracy of 73% was achieved on the Emergent test data set using this method. This method was challenged though in detecting observing stances, due to the similarities between the headlines and the claims when a headline is observing. This lead to a mislabelling of the observing claims as for claims.

Other work done on stance detection includes Stance Detection with Bidirectional Conditional Encoding (Augenstein et al., 2016). This paper looked at the stances of tweets, given sparse training data or where the target is not explicitly mentioned in the text. Two baselines were used, in a manner of treating stance detection as sentence level sentiment analysis. One was implemented using a Support Vector Machine Classifier and the other with long-short term memory networks (LSTM). LSTM networks were first proposed by Hochreiter and Schmidhuber and it is a gradient based method of learning to store information over extended time intervals (?). Their paper found LSTM networks to work best and were

used for most of the encodings.

Initially, the text and targets were independently encoded as a k-dimensional dense vector space, using two different LSTM networks. The model learned target-independent representation for the tweets, with the initial training encoding the tweets independant of any target. It relied on the nonlinear projection layer to incorporate the target in the stance prediction (Augenstein et al., 2016).

The paper also tested a conditional encoding method. First the target was encoded as a fixed length vector using one LSTM network. The tweet was then encoded using another LSTM network with its initial state using a representation of the target. This encoding was adapted to use bidirectional conditional encoding. One encoding was achieved by reading the target and the tweet from left-to-right, and then another by reading them from right-to-left. This allowed for the context on either side to be considered, in target dependant encoding method.

To deal with the small amount of training data, unsupervised pretraining was used, by initialized word embeddings used in the LSTM with a trained word vector model. These embeddings were only used for initialization and were optimized with further training. The paper found that conditional encoding was well suited for learning how to fit a targets with generalized encodings and that bidirectional encoding performed best overall, especially where the target was not explicitly mentioned in the tweet.

3 Deception Detection

Conroy et. al. discusses several approaches for automatic fake news detection using natural language processing. This is described as categorizing news on a spectrum based on their level of certainty as well as their veracity (the intention to mislead). The automatic detection of fake news is centered around predicting the chances that any news item is intentionally misleading based on the content of the news item. According to the research, the two main approaches currently being used are linguistic methods and network methods. Both of these methods make use of machine learning on their training data set (Conroy et al.,

200 2015).

204 3.1 Linguistic methods

206 These methods utilize knowledge of speech pat-
 207 terns that are able to identify truthfulness and de-
 208 ception more accurately than most humans. Under
 209 this approach, a basic way that text is analyzed is
 210 considering all words in a block of text as equally
 211 significant units. Using natural language process-
 212 ing, this technique would be implemented using
 213 *n-grams* to analyze word frequencies and find in-
 214 dicators of deception. This method might also in-
 215 volve tagging the lexical cues of words (also called
 216 shallow syntax) or frequencies of words which
 217 can uncover linguistic patterns of deception. The
 218 techniques under this approach rely heavily on the
 219 analysis of the usage of language. They also work
 220 very well when combined with other approaches
 221 (Conroy et al., 2015).

224 3.2 Network methods

226 These methods make use of a network of as-
 227 sociated information (like metadata) to predict
 228 the level of veracity of the content. It is also
 229 pointed out that the use of networks of data can
 230 provide a means to check the validity of a news
 231 item due to the presence of findable truths in
 232 the network. This involves making queries on
 233 existing knowledge to measure the truthfulness of
 234 new news items (Conroy et al., 2015).

236 The conclusion drawn was that both methods
 237 are very accurate in classifying news items. This
 238 gives rise to the use of a hybrid methods that takes
 239 into account both approaches to automating fake
 240 news detection. Such hybrid methods would have
 241 a linguistics-based analysis process that takes
 242 into account lexical analysis. It would also be
 243 able to perform efficiently in place of a strictly
 244 linguistic or network based approach. These
 245 techniques should be created with the intent of
 246 complementing the processes performed by a
 247 researcher in detecting fake news, as opposed to
 248 replacing them (Conroy et al., 2015).

250 4 Proposed Approach

251 4.1 Data Available

252 4.2 Baseline Implementation

253 4.3 Potential Improvements

254 255 256 References

- 257 Hunt Allcott and Matthew Gentzkow. 2017. Social me-
 258 dia and fake news in the 2016 election .
- 259 Isabelle Augenstein, Tim Rocktäschel, Andreas Vla-
 260 chos, and Kalina Bontcheva. 2016. Stance detec-
 261 tion with bidirectional conditional encoding. *arXiv*
 262 *preprint arXiv:1606.05464* .
- 263 Fake News Challenge. ????. Fake news chal-
 264 lenge stage 1 (fnc-i): Stance detection.
 265 <http://www.fakenewschallenge.org/>.
- 266 Niall J Conroy, Victoria L Rubin, and Yimin Chen.
 267 2015. Automatic deception detection: methods for
 268 finding fake news. *Proceedings of the Association*
 269 *for Information Science and Technology* 52(1):1–4.
- 270 William Ferreira and Andreas Vlachos. 2016. Emer-
 271 gent: a novel data-set for stance classification. In
 272 *Proceedings of the 2016 Conference of the North*
 273 *American Chapter of the Association for Computa-*
 274 *tional Linguistics: Human Language Technologies.*
 275 ACL.
- 276 Adam Kucharski. 2016. Post-truth: Study epidemiol-
 277 ogy of fake news. *Nature* 540(7634):525–525.
- 278 Jonathan Rose. 2017. Brexit, trump, and post-truth pol-
 279 itics.
- 280 Victoria L. Rubin, Yimin Chen, and Niall J. Con-
 281 roy. 2015. Deception detection for news: Three
 282 types of fakes. *Proceedings of the Association*
 283 *for Information Science and Technology* 52(1):1–4.
 284 <https://doi.org/10.1002/pr2.2015.145052010083>.
- 285 S Wineburg, S McGrew, J Breakstone, and T Ortega.
 286 2016. Evaluating information: The cornerstone of
 287 civic online reasoning.