

# Fake News Detection with Naive Bayes algorithm

\*Machine Learning Project

Muhammad Shahzeb Kayani  
Dept. of computer engineering  
Air University  
171209

**Abstract**—With the current political climate and parties more polarized than ever, the dissemination of accurate and unbiased information is more important than ever. However, there are limited resources for automatically identifying fake news articles on the Internet. Most modern machine learning pipelines that utilize deep neural networks (DNNs), which can model complex language features such as grammar, have been proven to be effective but computationally expensive.

**Index Terms**—Machine, Learning, Clustering, unsupervised learning

## I. METHODOLOGY

We treat the problem as a conditional probability task using Bayes' theorem, where our classifier attempts to predict the likelihood of a class (i.e. fake or real) given a document D. More specifically, we attempt to find the likelihood that a word W is likely to be from a fake news document with class F, and then combine those probabilities to generate a prediction for the whole document. The simple Bayes' formula can be expressed as follows.

$$P(F|W) = \frac{P(W|F)P(F)}{P(W)}$$

For our purposes, the extended form of Bayes' theorem was particularly useful, because  $P(W)$  could be expressed in terms of  $P(W|F)$  and  $P(W|\neg F)$ , under the assumption that classes F (fake) and  $\neg F$  (real) are mutually exclusive.  $P(W|F)$  is the frequency of word W in all fake documents in the training set,  $P(F)$  is the frequency of documents of class F.

$$P(F|W) = \frac{P(W|F)P(F)}{P(W|F)P(F) + (P(W|\neg F)P(\neg F))}$$

Finally, we calculate the likelihood of a document being fake, under the assumption that the likelihoods of words appearing in fake and real documents are independent, by multiplying the  $P(F|W_i)$  for all words in the given document. Because this is a binary classification problem, if the  $P(F|D)$  >  $P(\neg F|D)$ , we can infer that the document is likely to be fake.

## II. SOLUTION

The quantity of data required in order to produce accurate predictions is well beyond the realm of manual calculation, so

we wrote a program using the Python programming language to calculate the probabilities for the Bayes classifier. This also gave us the opportunity to further optimize our data set.

## III. DATASET

Our dataset came from a compilation of two datasets independently published on the popular opensource software website GitHub. The dataset was used in research published by Horne et. al, 2017. The first half of the dataset consisted of BuzzFeed Political News Data, categorized accordingly by analysts at BuzzFeed. The second half of the dataset included randomly collected political news data from real sources (i.e. Wall Street Journal, The Economist, BBC, NPR, USA Today, etc.) and fake sources (i.e. True Pundit, DC Gazette, Liberty Writers News, InfoWars, etc.). For the random political news partition, the sources were verified as real or fake based on Business Insider's "Most Trusted" list and Zimdars 2016 Fake News list.

For our purposes, we combined the two partitions of the original dataset into two categories, one sample of 128 real news articles and another sample of 123 fake news articles.

## IV. DOCUMENT PREPROCESSING AND TOKENIZATION

Using Regular Expression tokenization, we were able to remove punctuation and more easily segment words in a given document. We also removed stop words, or words that appear frequently and give little indication of class, such as "the", "and" or "as". Given stop words' frequency, even small variations in their usage could be enough to drown out the words that are good indicators of a particular class.

## V. RESULTS

Using the methodology described in Section 1, we trained our classifier on a sample of 128 real and 123 fake documents consisting of 128,858 words (71,731 after the removal of stop words). As a result, our program was able to predict the class of a document it had never seen before with 77.6 accuracy. The false positive rate for misclassifying a real news article as fake was 5.4 and the false negative rate was 38.5.

