

# Fake News Assignment

cmkv68

March 4, 2018

## Program Prerequisites

The python program utilises numpy, keras and nltk. Keras by default installs tensorflow, but it is recommended to install tensorflow manually to optimise its performance (as it's own installation compiles the program based on the user's machine, whereas installing via pip3 does not.) The program will attempt to download the GloVe dataset itself.

```
pip3 install numpy nltk tensorflow keras
python3 main.py
```

## Testing Configuration

Both learning methods are subject to a machine utilising an `intel i7-6600u` with 12 Gigabytes of DD4 memory, alongside NVMe-E storage. Due to the relatively large nature of the word2vec dataset which will be used, running times may vary dependent the hardware utilised such as the storage, memory and the processor itself.

## Data Sanitisation

Each document's article is converted into lowercase and is then subject to the ordered removal of:

- Strong quotation marks
- Commas between numbers (and joining numbers together)
- Apostrophes between letters, alongside the suffix (which indicate plurality, possessive cases or contractions)
- URLs
- Line breaks
- Unicode strings
- Duplicate quotes
- Punctuation
- Stopwords

The removal is ordered as removing punctuation entirely in one farce may produce unintended effects, i.e `20,000,000 => 20 000 000` or `John Lennon's => john lennon s`.

The reasoning behind the removal of stopwords is related to the fact that the words themselves are used for grammatical understanding, which are not considered in our learning methods, especially in the case for shallow learning.

Once computed, a dataset object is created, containing each document in the dataset. Each document is an object itself, containing the raw article, and the cleaned data, (now a list of words in python). While this is

not necessarily the most efficient method, it was chosen to help understand the learning mechanisms used in the assignment.

## Shallow Learning

The majority of the shallow learning was implemented manually to help understand the inner workings of how shallow learning techniques work.

### Term Frequency

Each document is iterated, with each word being counted relative to the rest of the other words in the document. This is then stored in the document's object in a dictionary, where each word is a key, and their word count is the value.

### Term-Frequency/Inverse Document Frequency

A global document frequency  $\mu$  for words is introduced, using the training documents's words. For each document, we iterate through its term frequency and accumulate them to  $\mu$ . This ensures that each unique word is counted per document as opposed to their word-count/term-frequency. Once this is accumulated, Each document is iterated through and a TF-IDF value is calculated for each word  $w$  using the formula  $tfidf(w) = tf(w) * df(w)^{-1}$ . TF-IDF values are stored in the same style as how term-frequency is stored in a document's object.

### N-Grams

There is the option to use either bigrams or trigrams (unigrams can be considered equal as term-frequency.), when using the n-grams option. Using the NLTK library, N-Grams are produced, using the document's cleaned word-list and are stored in the document's object.

```
trigrams([the quick brown fox jumped over the lazy dog])

(the quick brown)
(quick brown fox)
(brown fox jumped)
...
(the lazy dog)
```

### Naive Bayes

Words that are found in the test document but not in the training document are **ignored**. This is chosen to reinforce the worth of the training data.

With  $\alpha$  as the conditional probability of a word or a wordgroup for a given class,  $\epsilon$  being the class itself, the naive bayes classifier is composed of the formula:

$$\beta = (count(\alpha, \epsilon) + 1) / (count(\epsilon) + vocabulary)$$
$$P(X|Y) = \beta \cdot \gamma$$

Where  $\gamma$  represents the probability of a class.

This formula is performed for each of the classes (fake and real), and the class is chosen based on the maximum value produced between the classes.

## Deep Learning

For deep learning, we subject each document using a public dataset, the glove word2vec model (which can be found at <https://nlp.stanford.edu/projects/glove/>) as a premise to compare the documents against.

There are t

## Results

Additional Assessment Criteria: A.General Performance of the solution on the test data set -Are the results comparable or above the expected baseline (i.e. > 75% accuracy)? (10 Marks) -How all the components work together to achieve the reported results (10 Marks)

### Shallow

Method	Precision	Accuracy	Recall	F1 Measure	Time Taken
TF	49.81	87.68	87.92	63.59	Negligible
TF-IDF	94.19	51.67	50.87	66.06	Negligible
Trigrams	50.34	99.33	98.68	66.67	Negligible
Bigrams	50.34	99.33	98.68	66.67	Negligible

Out of our shallow learning methods, we can see

TF-IDF performs the worst in terms of accuracy, which makes sense. Words that are necessarily unique to a given test document may not be considered on training data, which makes it difficult to precisely quantify the significance of those particular words. It also happens to be the case that those words would hold more weight than other words that are considered otherwise in other documents.

In other terms, TF-IDF doesn't differentiate between the terms and the perspective of conveyed sentiments (as is the case with all shallow learning methods). For instance, we assume the training set has the words "toaster", "bad", "good" and the test data contains "magnet", "good". Since we have a weighting for "good" as it is in the training set, we do not have any information on the term "magnet", even though the term itself holds more semantic information than "good".

This may also help show why the precision score of TF-IDF is especially high compared to the other methods.

### Deep

Method	Precision	Accuracy	Recall	F1 Measure	Time Taken
LSTM (E1)	64.0	85.667	73.909	67.234	40s
CNN (E1)	55.376	66.0	54.152	50.386	23s
LSTM (E2)	74.667	90.333	61.606	67.331	41s
CNN (E2)	62.359	69.333	32.303	42.474	23s

Running more epochs can naturally improve the quality of the classifier with the tradeoff of introducing possible overfitting of the data.

Due to the inherent differences between the two approaches, it makes it very difficu

## **Conclusion**