

Fake News Detection

Big Data for Official Statistics
February 2021

Oluwatoyin Yetunde SANNI
Matricola: 1871835



What is fake news?

Fake news—news articles that are intentionally and verifiably false designed to manipulate people's perceptions of reality—has been used to influence politics and promote advertising. But it has also become a method to stir up and intensify social conflict.

Why should we care if a news is fake?

It has been used to influence politics and promote advertising. But it has also become a method to stir up and intensify social conflict.

It is now known that false information played a major role in the last American presidential election.

Michelle Obama Deletes Hillary Clinton From Twitter

When Hillary goes low, Michelle goes BYE!

Posted on November 1, 2016 by Baxter Dmitry in News, US // 43 Comments



FALSE

Fake news—news articles that are intentionally and verifiably false designed to manipulate people's perceptions of reality—has been used to influence politics and promote advertising.

Fake news—news articles that are intentionally and verifiably false designed to manipulate people's perceptions of reality—has been used to influence politics and promote advertising.



Cindy Otis (Pre-order TRUE OR FALSE now!)

@CindyOtis_



1. THREAD.

There's a fake news story circulating from a website called "MCM News" claiming the Pope has the [#coronavirus](#). The domain was registered in 2016 by a domain squatter in China. The registration was changed three days ago on 26 Feb.

📅 Saturday, February 29 2020 🔒 Marzy Star co-founder David Raback dies, aged 61



MCM News



👍 643 10:11 AM - Feb 29, 2020



Impacts of Fake news spreads across different sectors around the world

A decorative network diagram in the top right corner, featuring a series of interconnected nodes (circles) of varying sizes, some solid and some hollow, connected by thin lines, suggesting a social or information network.

Democratic impacts

Will I vote differently if I know the Pope endorses a politician's candidacy?

Fear

Will I have insomnia?
Will my quality of life decrease?

Health Impacts

Will I take the Covid19 Vaccine if I read a news article that it killed or has negative lasting effects?

Financial impacts

If I like that blogger and her post pushes me to buy a certain product?

A decorative network diagram in the bottom left corner, similar to the one in the top right, showing a cluster of interconnected nodes and lines.

TOP
5 FAKE
NEWS

During
trump
election

Fake News Is A Real Problem

Facebook engagement of the top five fake election stories*



Total Facebook engagement for top 20 election stories (August-election day)



Where Exposure To Fake News Is Highest

% who say they were exposed to completely made-up news in the past week*



TOP
5 FAKE
NEWS

During
trump
election

* Selected countries
n=74,000 respondents in 37 markets (Jan/Feb 2018)

Source: Reuters Institute Digital News Report 2018

The major goal of what I have done in this project is to classify news content as Fake/Real

Hence, I did the following in this project:

1. Obtain the datasets.
2. Preprocess the data and prepare it in a format that machine learning models can work with.
3. Did some Data exploratory analysis to understand my Data.
4. Built different machine learning and deep learning models to learn the prepared data to compare results and achieve the best possible classification results.





“

Dataset and Preprocessing

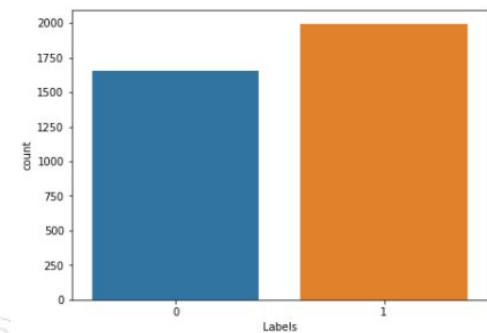
“As data scientists, our job is to extract signal from noise.” —Daniel Tunkelang

Dataset

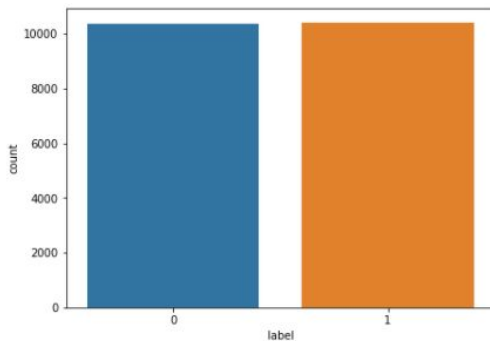
- ◎ I combined two datasets from [MachineHack](#) and Kaggle. They both consists of texts and the corresponding binary classification.
- ◎ Since this is a news article, the datasets isn't heavily noisy. However, it is still extremely difficult to learn the specific words and sequences that matter to detect a false news content. This also makes the vocabulary bloated , and the tweet vectors sparse.

Dataset

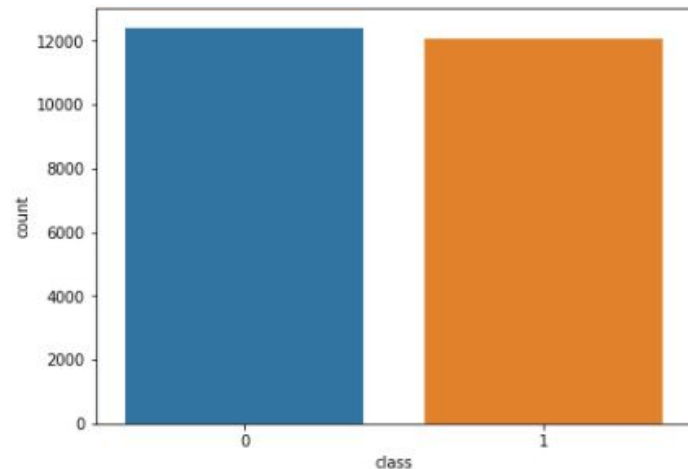
Interestingly, both have a normal distribution with an even distribution of the target. This gives us a balanced dataset which was splitted into 70% for training and 30% for cross validation.



+



=



Data Cleaning and Preprocessing

- ◎ Since our task takes consideration of context and semantic meanings, I parsed the dataset to remove words that will provide no intuition. Words like punctuations, stop words, were cleaned off as well.
- ◎ Words with accents were represented in Unicode from the dataset, so I removed them.
- ◎ I also converted all the words to lowercase, also to reduce the vocabulary size.

Data Cleaning and Preprocessing

Our models can only understand numbers, hence representing our texts in numeric is a crucial step. This is called feature extraction or vectorization. For this project, I explored Count Vectorizer and TF-IDF(Term Frequency Inverse Document Frequency).

Count Vectorizer

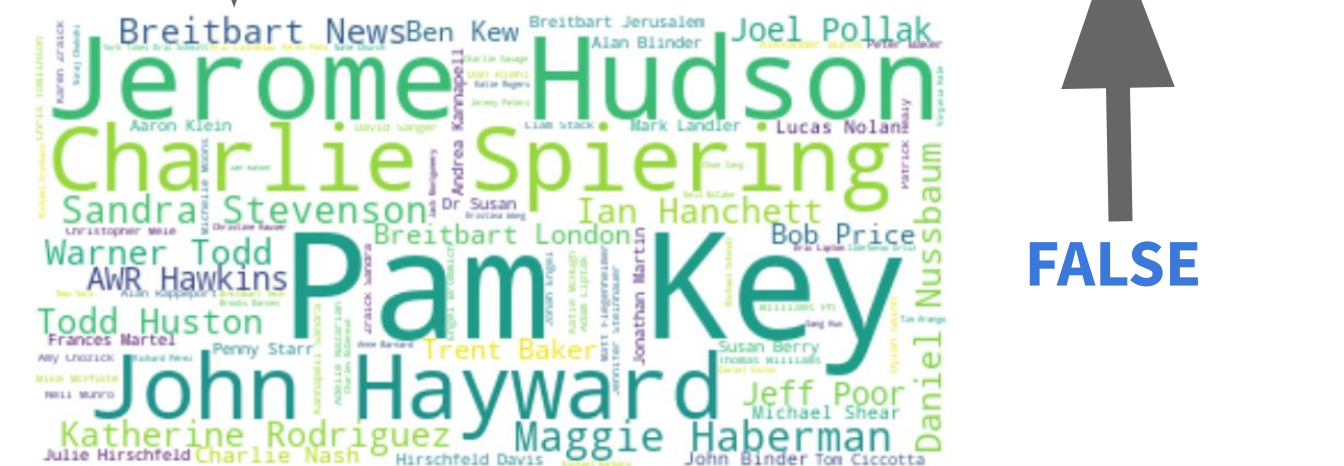
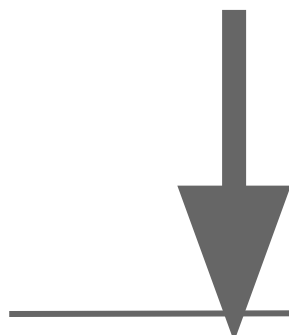
	blue	bright	sky	sun
Doc1	1	0	1	0
Doc2	0	1	0	1

TD-IDF Vectorizer

	blue	bright	sky	sun
Doc1	0.707107	0.000000	0.707107	0.000000
Doc2	0.000000	0.707107	0.000000	0.707107

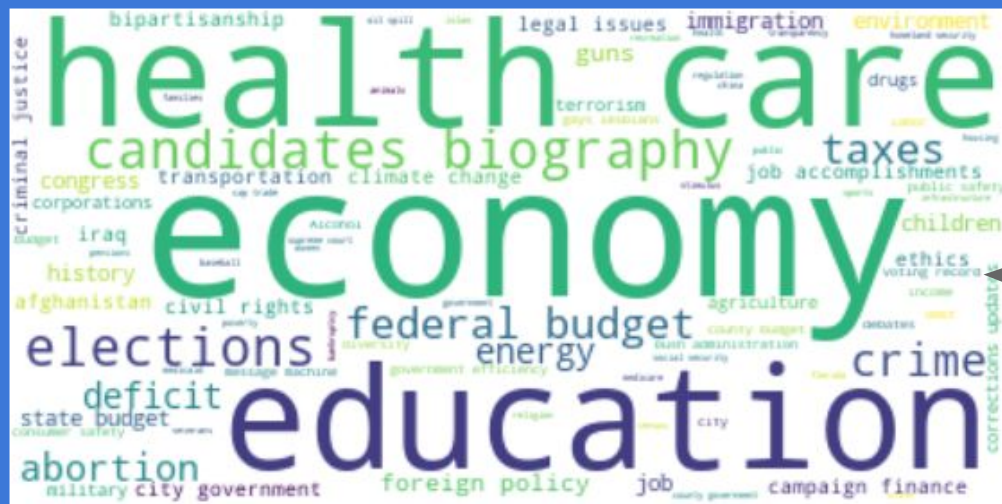
Word Count of the most frequent **Authors** for False and True News content

TRUE

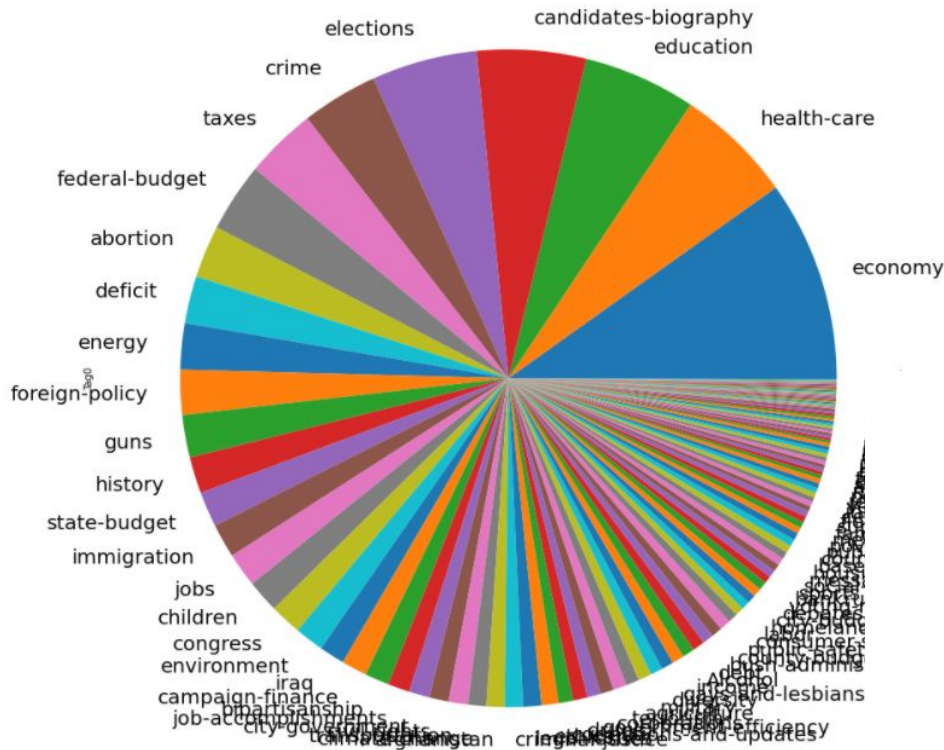


FALSE

FALSE

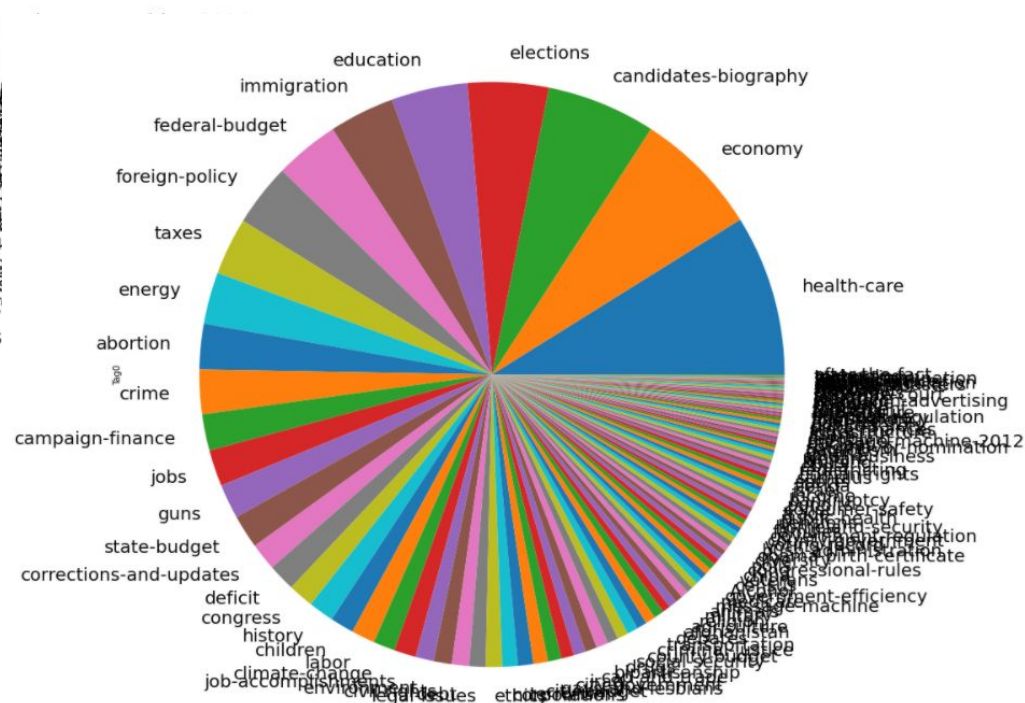


The Pie Chart gives us more insights to our previous Word Count



TRUE

FALSE





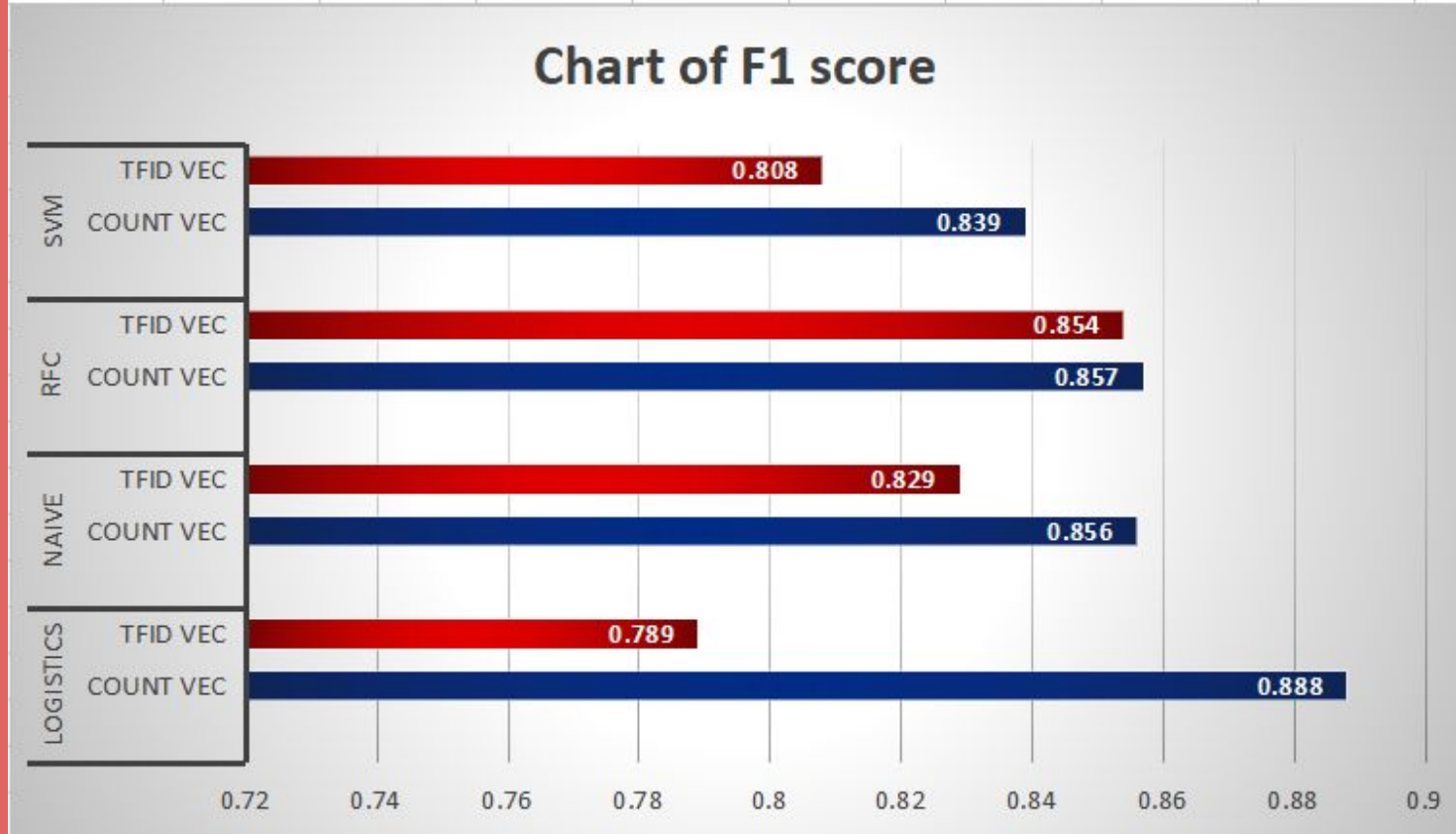
“

Experiments and Results

*“All knowledge - past, present, and future
- can be derived from data by a single,
universal learning algorithm.” — Pedro
Domingos*

Experiment One

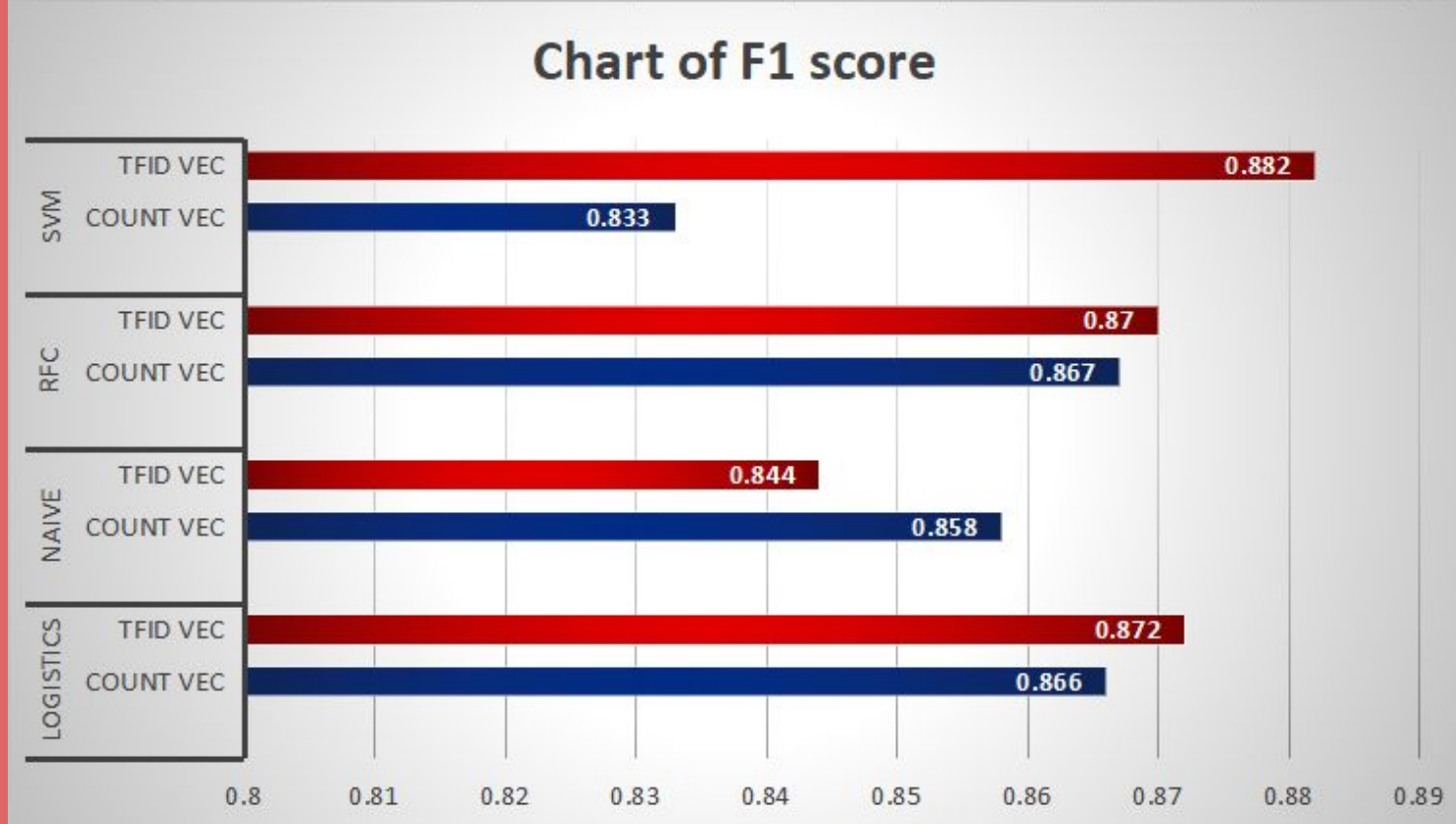
Selecting the best Statistical ML model among 4 classifiers.



The chart above shows the F1 scores on test data **without applying the data cleaning process**. Logistics Regression with count vector features features better generalization of the data.

Experiment Two

Selecting the best Statistical ML models among 4 classifiers.

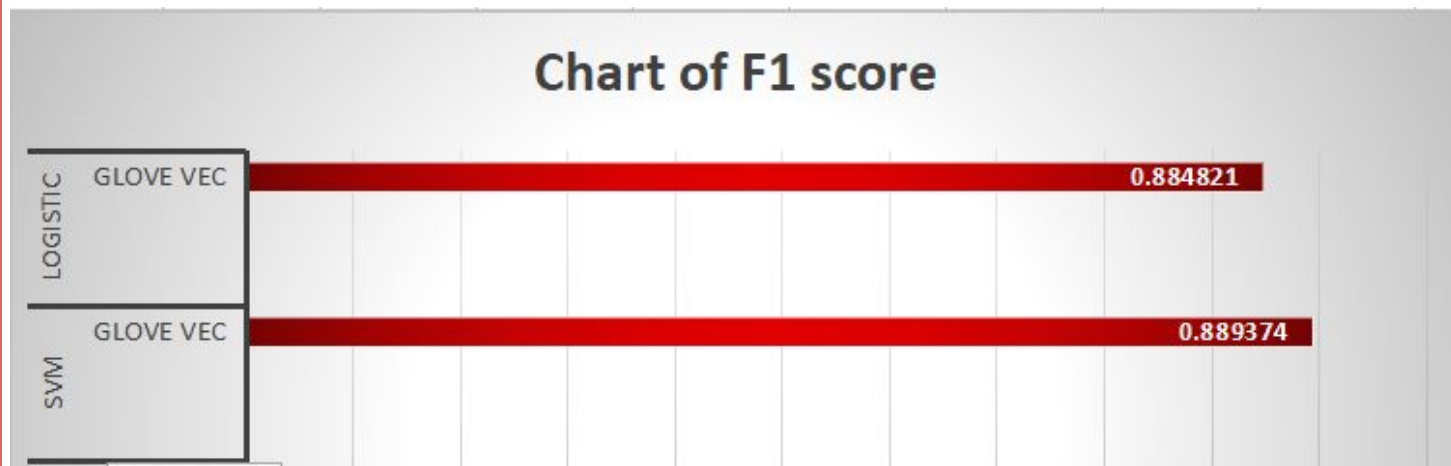


The chart above shows the F1 scores on test data from this experiment. **SVM with TF IDF** weighted features shows better generalization of the data.

Experiment Three

Selecting the best Statistical ML model among 4 classifiers.

I also combined the use of Glove Embeddings with SVM and Logistics Regression, this experiment shows a slight improvement in the dataset.



The chart above showing the F1 scores on test data **reveals SVM** performed better.

Selecting the best Statistical ML model

Model	SVM	Logistics
Cleaned Data	No	Yes
Embedding used	TF IDF	Count Vectorizer
F1 Score	0.882	0.888

I conducted a parameter grid search on the SVM model was able to improve the F1 score on dataset 2 to 0.735 with the TF IDF vectorizer parameters `max_df = 0.5` and `ngram_range = (1, 2)`).

SVM: 0.904	LOGISTICS: 0.882
<code>max_df = 0.5,</code> <code>ngram_range = (1, 2)</code> <code>C value = 1.</code>	<code>max_df = 0.5,</code> <code>ngram_range = (1, 2)</code> <code>C value = 0.1</code>

Experiment Four

Can Neural
Networks do
better using
BiLSTMs?

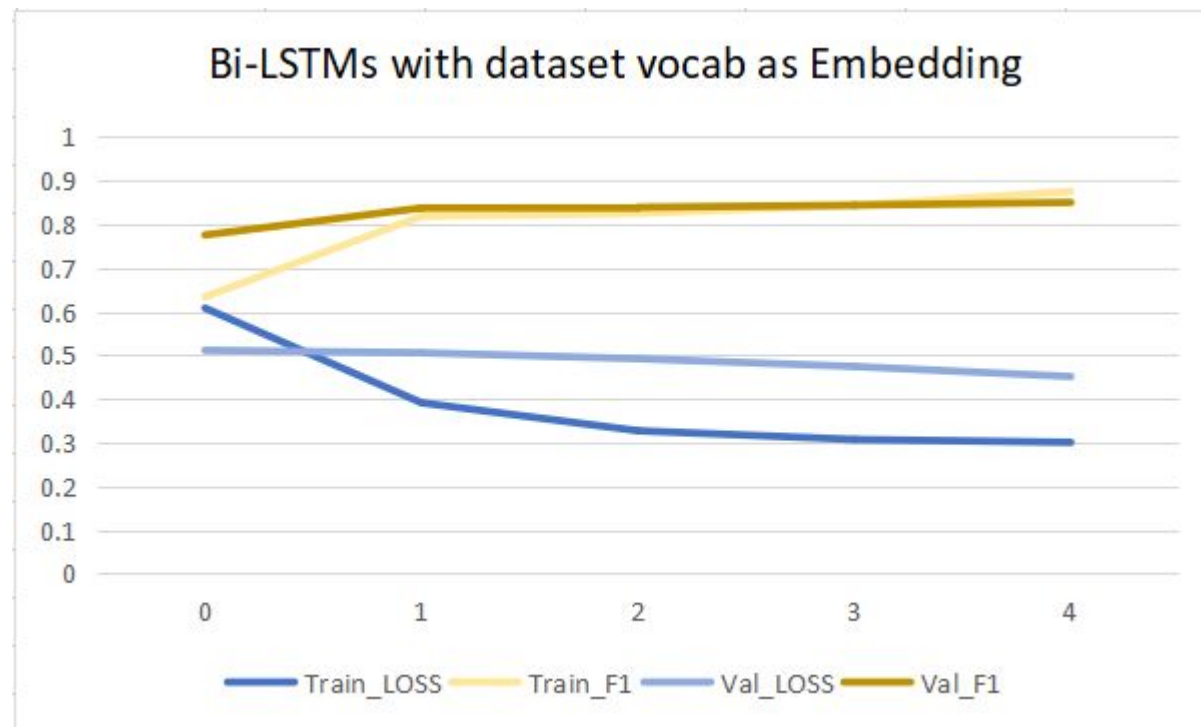
Bi LSTMs have performed really well for various NLP tasks because of their ability to learn word sequences, so I explored two Bi-LSTM model configurations, and a transfer learning on a pre trained model.

1. Bi LSTMs with our dataset vocabulary as the **embedding layer**.
2. Bi LSTMs with **GloVe pre trained embedding** word Vectors.

Experiment Four

Can Neural Networks do better using BiLSTMs?

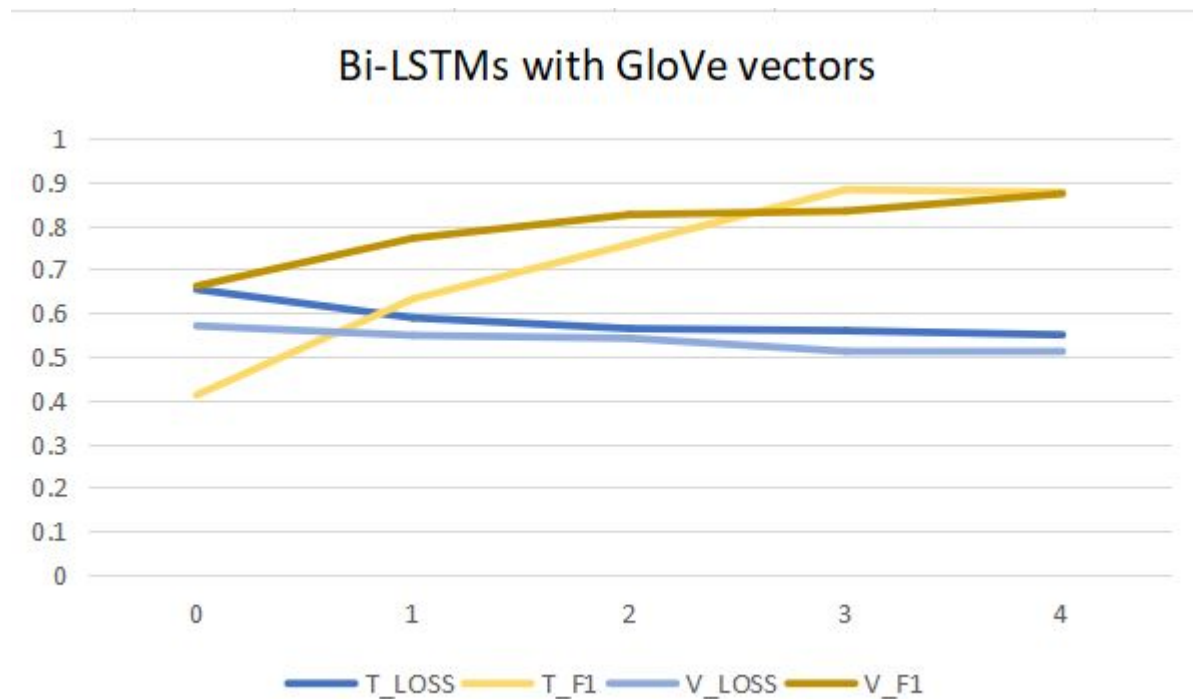
Bi-LSTMs with dataset vocab as Embedding	Bi-LSTMs with GloVe vectors
0.936	0.947



Experiment Four

Can Neural Networks do better using BiLSTMs?

Bi-LSTMs with dataset vocab as Embedding	Bi-LSTMs with GloVe vectors
0.936	0.947



CONCLUSION FROM EXPERIMENTS

1. Statistical ML models performs well on sparse data representations while deep learning approaches on sparse data representations because they have no information of context semantics and sequence.
2. RNNs perform greatly with sequential data especially when the data size is large. Hence, combining our dataset was a good call.
3. Word embeddings are good for representing semantics and meaning, this is a fact why the deep learning models do well with this representation.



A decorative background featuring a network diagram. It consists of numerous nodes, represented by small circles, connected by thin lines. Some nodes are solid blue, while others are grey with a blue outline. The network is more densely packed on the left and right sides of the image, with the central area being mostly white space containing the text.

**THANKS
FOR
LISTENING**

References

1. <https://www.webwise.ie/teachers/what-is-fake-news/>
2. <https://www.cits.ucsb.edu/fake-news/danger-social>
3. [https://monkeylearn.com/text-classification/#:~:text=Some%20of%20the%20most%20popular,SVM\)%2C%20and%20deep%20learning.](https://monkeylearn.com/text-classification/#:~:text=Some%20of%20the%20most%20popular,SVM)%2C%20and%20deep%20learning.)
4. <https://medium.com/the-innovation/sentiment-analysis-using-word2vec-and-glove-embeddings-5ad7d50ddb0d>
5. <http://nadbordrozd.github.io/blog/2016/05/20/text-classification-with-word2vec/>

