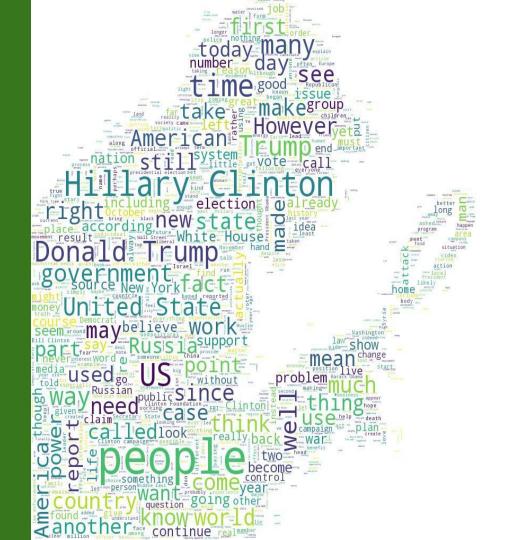# FAKE NEWS CLASSIFIER

## Data Crunchers

Tejas Prashanth, Suraj Aralihalli,
Sumedh Basarkod, Sumanth Rao

1. Recognizing style of writing
2. Fact-checking

# Module 1:
# Writing style recognition

# Features Considered

1. Term Frequency ( TF-IDF )
2. Readability
3. Punctuation Usage
4. Polarity

# Term-Frequency

- Identifying the style of writing by using the frequency of words
- Based on how important a word is to a document

# Readability

9 Metrics that score understandability of text which are based
9 different formulae which include

- Linear Write Index
- Dale-Chall Readability Score
- Automated Readability Index
- SMOG Index
- Flesch Reading Index
- Flesch Kincaid Grade
- Linear Write Index
- Difficult Words

| | |
|---|---|
| 4.9 or lower | average 4th-grade student or lower |
| 5.0 to 5.9 | average 5th-grade or 6th-grade student |
| 6.0 to 6.9 | average 7th-grade or 8th-grade student |
| 7.0 to 7.9 | average 9th-grade or 10th-grade student |
| 8.0 to 8.9 | average 11th-grade or 12th-grade student |
| 9.0 to 9.9 | average 13th-grade to 15th-grade student |

# Punctuation Usage

- Usage of upper-case alphabets
- Normalized count of punctuations which include ! ? - , . $ etc
- Normalized Word Count
- Count of blank spaces ( " " ) normalized by count of full stops ( . )

# Polarity

Example:

Like you, I am getting very frustrated with this process. I am genuinely trying to be as reasonable as possible. I am not trying to "hold up" the deal at the last minute. I'm afraid that I am being asked to take a fairly large leap of faith after this company has screwed me and the people who work for me.

- Negative : 0.093
- Neutral : 0.836
- Positive : 0.071

# Hybrid Model 1

- Features are:
  1. Readability features
  2. Punctuation usage
  3. Polarity

- 22 features in total

# Classifiers and Accuracy

| Model | Accuracy |
|---|---|
| Logistic Regression | 61 % |
| AdaBoost Classifier | 61 % |
| Random forest | 64.3 % |

# Hybrid Model 2

- Features include
  1. Readability
  2. Punctuation usage
  3. Polarity
  4. TF-IDF

- We chose this model over the previous model as TF-IDF could give promising results when combined with other engineered features

- Logistic regression with Stratified 3-fold validation gave an accuracy of 64 %

# Module 2:
# Fact-checking

# Preprocessing

1. Conversion to lowercase
2. Keyword extraction
3. Lemmatization

# Fact extraction approach :

1. Semantic and dependency parsing is performed
2. Sentences with named entities in subject, object are selected :

   ((subject, object, action), True/ False)

3. Technology used: IBM Watson Semantic Roles extraction

```
(('artificial intelligence AI Advances machine learning ', 'workforce ', 'have'), False)
(('artificial intelligence AI Advances machine learning ', 'workforce ', 'impacted'), False)
(('humans relationship machines ', 'human skills ', 'will likely continue to augment'), False)
(('report ', 'them', 'than replace'), False)
(('report ', 'artificial intelligence AI human skills workforce Advances ', 'noted'), False)
(('artificial intelligence AI Advances machine learning ', 'workforce ', 'have'), False)
(('artificial intelligence AI Advances machine learning ', 'workforce ', 'impacted'), False)
(('humans relationship machines ', 'human skills ', 'will likely continue to augment'), False)
(('report ', 'them', 'than replace'), False)
(('report ', 'artificial intelligence AI human skills workforce Advances ', 'noted'), False)
(('PwC consultancy firm report ', 'artificial intelligence breakthroughs threat AI ', 'found'), False)
(('NEWS ', 'latest fintech news globe ', 'Contents'), False)
(('by human', 'certain jobs ', 'performed'), False)
```

# Fact validation approach

a. Search API phase - Top 10 search results are fetched for the fact that is to be validated
b. Semantic and dependency parsing- Identify subject, object and actions of each search result
c. Keyword extraction-Extract keywords from text of subject, object and action

# Fact validation approach

d.    Lemmatization-Convert verbs to their root forms

e.    Document embeddings-Extract word embeddings of each word and find average value

f.    Cosine similarity-Compare the embeddings of two documents using cosine similarity.

i.    High similarity: Accurate fact
ii.   Low similarity: Inaccurate fact

# Sample output

```
sumanthvrao@sumanthvrao-inspiron-7460:/media/sumanthvrao/Personal/PESU/OTHERS/PRAVEGA_2019/Pravega_2019_FINAL_ROUND/Data_C
runchers$ ./classifier ./training/fakeNewsDataset/legit/entmt01.legit.txt
Reid
----------------
Joseph
----------------
Marcell
----------------
Geoffrey
----------------
1.0
1.0
1.0
----------------
MATCHED subject Karyn Parsons Hilary Tatyana Ali Ashley Daphne Maxwell Reid Joseph Marcell Geoffrey   :::   Karyn Parsons H
ilary Tatyana Ali Ashley Daphne Maxwell Reid Joseph Marcell Geoffrey
MATCHED Object Ribeiro Smith   :::   Ribeiro Smith
MATCHED Action were joined  :::  were joined
===================================
The fact mentioned above is True.
Accuracy of predicting Fact is True :  1.0
/home/sumanthvrao/anaconda3/lib/python3.6/site-packages/sklearn/base.py:311: UserWarning: Trying to unpickle estimator Log
isticRegression from version 0.19.0 when using version 0.19.1. This might lead to breaking code or invalid results. Use at
 your own risk.
  UserWarning)
Article is legitimate according to classifier
```

# Combination

- We analyse both the models and take a hybrid combination of both.
- A combination of scores from the Fact-checker and from the writing style recognition model is used.

# Future Work

- Improving the accuracy Fact-checker
  - Using Relation extraction for extracting facts
- Using topic classification
  - Applying hybrid model to each topic separately