**Fake News Classifier**:
Pinky Chauhan
University of Illinois at Urbana Champaign

## Overview:

- This objective of this project is to build a classifier system based on machine learning that is able to identify fake news from real/reliable news given a news title and/or news text content as the input. Such a tool can be integrated with social media platforms to flag potentially fake articles or filter those out.
- This is essentially a data categorization problem where I have trained several classifier models on the following dataset from Kaggle: https://www.kaggle.com/c/fake-news/data
- The dataset comprises of csv format files for training and testing with each file containing id, news title, text and author fields.
- The train.csv also has label field to categorize data as Reliable (label value 0) and Fake (label value 1)
- After evaluation based on various performance metrics, one of the models (in this case, Linear SVC over unigram bag-of-words/TF-IDF representation) is integrated in the final tester notebook to test with news data.
- The classifier takes a news article (title and text) as input and provides a prediction for the news article as either of the 2 categories:
  - Fake News
  - Reliable News

## Software Implementation Details:

• Data Analysis:
  - Comprises of checking several attributes to evaluate their contribution towards classification. For such a classifier, the text and title of the news make obvious choices as features.
  - I also analyzed authors' distribution using pandas and polarity/sentiment differences using NLTK vader sentiment intensity analyzer library on the dataset.

• Preprocessing of data:
  - Handling missing values by removing any rows with no text and title, preprocess data to remove any punctuations, remove any words with length 3 or less, stop words removal, tokenization and lemmatization using NLTK libraries

• Feature selection:
  - Concatenated news title and text into article field and preprocessed it. Article comprises the feature to train the model

• Vectorization

- Different vector forms listed below have been used using NLTK vectorization/transformation libraries:
    - Term frequency (TF) based vector over unigrams bag of words representation
    - Term frequency/inverse document frequency (TF-IDF) based vector over unigrams
    - Term frequency (TF) based vector over unigrams and bigrams
    - Term frequency/inverse document frequency (TF-IDF) based vector over unigrams and bigrams
    - Term frequency/inverse document frequency (TF-IDF) based vector over unigrams, bigrams and trigrams

- Training/hyperparameter tuning/validation using classification models:
    - Models used (sklearn libraries):
        - Naïve bayes (With/without smoothing, TF vs TF-IDF vectors, Unigram/N-gram)
        - Logistic Regression (TF-IDF vectors using Unigrams/N-grams)
        - SVM using Linear SVC (TF-IDF vectors using Unigrams/N-grams, Regularization)
        - SGDC classifier (TF-IDF vectors using Unigrams/N-grams)
        - Decision Tree (TF-IDF vectors using Unigrams/N-grams)

- Performance evaluation:
    - Compute and analyze metrics using sklearn metrics libraries
        - Precision (macro/micro), recall (macro/micro), F1 (macro/micro)
        - Classification Accuracy
        - Confusion matrix to see distribution of true/false positives/negatives
    - Select the best performing model based on evaluation results (SVM using Linear SVC using TF-IDF vector over unigrams)
    - Results:

| | accuracy | precision(macro) | precision(micro) | recall(macro) | recall(micro) | f1_score(macro) | f1_score(micro) |
|---|---|---|---|---|---|---|---|
| Decision Tree (TFIDF/Uni-bi-trigram) | 96.13% | 0.961 | 0.961 | 0.961 | 0.961 | 0.961 | 0.961 |
| Decision Tree (TFIDF/Uni-bigram) | 95.13% | 0.952 | 0.951 | 0.951 | 0.951 | 0.951 | 0.951 |
| Decision Tree (TFIDF/Unigram) | 92.98% | 0.93 | 0.93 | 0.93 | 0.93 | 0.93 | 0.93 |
| Linear SVC (TFIDF/Uni-bi-trigram) | 95.60% | 0.956 | 0.956 | 0.956 | 0.956 | 0.956 | 0.956 |
| Linear SVC (TFIDF/Uni-bigram) | 96.21% | 0.962 | 0.962 | 0.962 | 0.962 | 0.962 | 0.962 |
| Linear SVC (TFIDF/Unigram) | 96.12% | 0.961 | 0.961 | 0.961 | 0.961 | 0.961 | 0.961 |
| Linear SVC (TFIDF/Unigram/Regularization) | 87.90% | 0.88 | 0.879 | 0.879 | 0.879 | 0.879 | 0.879 |
| Logistic Regression (TFIDF/Uni-bi-trigram) | 93.50% | 0.935 | 0.935 | 0.935 | 0.935 | 0.935 | 0.935 |
| Logistic Regression (TFIDF/Uni-bigram) | 94.00% | 0.94 | 0.94 | 0.94 | 0.94 | 0.94 | 0.94 |
| Logistic Regression (TFIDF/Unigram) | 94.54% | 0.945 | 0.945 | 0.945 | 0.945 | 0.945 | 0.945 |
| Multinomial naive bayes (TF/Uni-bigram/Smoothing) | 89.29% | 0.908 | 0.893 | 0.894 | 0.893 | 0.892 | 0.893 |
| Multinomial naive bayes (TF/Unigram/NoSmoothing) | 91.67% | 0.92 | 0.917 | 0.917 | 0.917 | 0.917 | 0.917 |
| Multinomial naive bayes (TF/Unigram/Smoothing) | 89.29% | 0.902 | 0.893 | 0.893 | 0.893 | 0.892 | 0.893 |
| Multinomial naive bayes (TFIDF/Uni-bi-trigram/Smoothing) | 82.15% | 0.867 | 0.822 | 0.821 | 0.822 | 0.816 | 0.822 |
| Multinomial naive bayes (TFIDF/Uni-bigram/Smoothing) | 76.87% | 0.839 | 0.769 | 0.771 | 0.769 | 0.757 | 0.769 |
| Multinomial naive bayes (TFIDF/Unigram/NoSmoothing) | 91.79% | 0.92 | 0.918 | 0.917 | 0.918 | 0.918 | 0.918 |
| Multinomial naive bayes (TFIDF/Unigram/Smoothing) | 82.69% | 0.869 | 0.827 | 0.825 | 0.827 | 0.821 | 0.827 |
| SGDC (TFIDF/Uni-bi-trigram) | 95.17% | 0.952 | 0.952 | 0.952 | 0.952 | 0.952 | 0.952 |
| SGDC (TFIDF/Uni-bigram) | 95.90% | 0.959 | 0.959 | 0.959 | 0.959 | 0.959 | 0.959 |
| SGDC (TFIDF/Unigram) | 95.75% | 0.957 | 0.958 | 0.958 | 0.958 | 0.957 | 0.958 |

- Save/export trained model:

- Using pipeline to specify all steps (vectorizer/classifier), fit training data and exporting model using joblib library

- Kaggle submission:
  - Predicted results for data in test.csv and submitted notebook/results to Kaggle (https://www.kaggle.com/pinkychauhan/fakenewsclassifierusingnltk-sklearn)
  - Accuracy: 94%

- Create script (Jupyter notebook) that will take news text as input and generate classification as reliable news or fake news.


**Installation/Execution Details**:
Code is written using Jupyter notebook and python 3

Code structure:
- data: This directory contains the dataset from Kaggle (https://www.kaggle.com/c/fake-news/data). There are 3 files:
  - train.csv: To use for analysis, training, validation
  - test.csv: Test dataset for submission of results to Kaggle competition
  - submit.csv: File containing results/predictions for data in test.csv
- notebooks: This directory contains 2 notebooks:
  - FakeNewsClassifierTraining.ipynb: Jupyter notebook containing code/results for data analysis, cleanup, features set up, vectorization, training using various classifier algorithms, tuning and performance evaluation/comparison, model pipeline creation/export, prediction of results for test.csv for Kaggle submission
  - Tester.ipynb: This notebook loads the pretrained/exported model and predicts the category for a given news article. Use this notebook to test the classifier.
- model: This directory contains the pretrained model exported by FakeNewsClassifierTraining.ipynb notebook and loaded by Tester.ipynb
- results: This directory contains the summarized performance metrics from different models used for training and a copy of the submit.csv file generated from predictions for data in data/test.csv

Code Setup:
- Install python 3 and Jupyter notebook
- Install the following python/machine learning libraries:
  - re: For regular expression matching
  - itertools: To iterate over data
  - pandas: For Data analysis/representation as Dataframes
  - nltk: Natural language toolkit
  - sklearn: For model selection, training, evaluation, export using pipeline
  - matplotlib: For visualization
  - joblib: For model export and load

- Checkout the project from main branch in Github
- Launch Jupyter notebook and navigate to the directory where project is checked out
- Tester.ipynb located in notebooks folder can be used for testing the classifier by providing values for title and text
- FakeNewsClassificationTraining.ipynb can also be executed to see all stages entailed in bulding the classifier and training/evaluation of different models

Note: In case you see an issue around missing packages stopwords, punkt, vader_lexicon or wordnet, download them one time using below commands:
nltk.download('vader_lexicon')
nltk.download('punkt')
nltk.download('stopwords')
nltk.download('wordnet')


**References**:
- https://www.kaggle.com/c/fake-news/data
- https://scikit-learn.org/stable/user_guide.html
- https://medium.com/datadriveninvestor/python-data-science-getting-started-tutorial-nltk-2d8842fedfdd
- https://matplotlib.org/tutorials/introductory/pyplot.html
- https://towardsdatascience.com/machine-learning-classifiers-a5cc4e1b0623