# Fake news detection using NLP

- **Introduction**

We consume news through several mediums throughout the day in our daily routine, but sometimes it becomes difficult to decide which one is fake and which one is authentic.

Do you trust all the news you consume from online media?

Every news that we consume is not real. If you listen to fake news it means you are collecting the wrong information from the world which can affect society because a person's views or thoughts can change after consuming fake news which the user perceives to be true.

Since all the news we encounter in our day-to-day life is not authentic, how do we categorize if the news is fake or real?

In this article, we will focus on text-based news and try to build a model that will help us to identify if a piece of given news is fake or real.

- **Fake News**

A sort of sensationalist reporting, counterfeit news embodies bits of information that might be lies and is, for the most part, spread through web-based media and other online media.

This is regularly done to further or force certain kinds of thoughts or for false promotion of products and is frequently accomplished with political plans.

Such news things may contain bogus and additionally misrepresented cases and may wind up being virtualized by calculations, and clients may wind up in a channel bubble.

- **Feature engineering in fake news detection**

Feature engineering is a critical step in fake news detection using Natural Language Processing (NLP). It involves transforming the raw text data into meaningful numerical features that can be used as input for machine learning models. Here are some common feature engineering techniques for fake news detection:

- Term Frequency-Inverse Document Frequency
- Word Embeddings
- Bag of Words (BoW)
- N-grams
- Sentiment Analysis

- **Model training evaluation in fake news detection**

Model training and evaluation are crucial steps in the development of a fake news detection system. Here's how you can approach them:

- Data Splitting

- Feature Engineering

- Model Selection

- Hyperparameter Tuning

- Training

- **Program**

```python
# Import necessary libraries
import pandas as pd

from sklearn.model_selection import train_test_split

from sklearn.feature_extraction.text import TfidfVectorizer

from sklearn.naive_bayes import MultinomialNB

from sklearn.metrics import accuracy_score, classification_report


# Load your labeled dataset
data = pd.read_csv('fake_news_dataset.csv')  # Replace with your dataset


# Split the dataset into training and testing sets
X = data['text']  # 'text' is the column containing news articles

y = data['label']  # 'label' is the column containing labels (0 for real, 1 for fake)

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)


# Feature Engineering - TF-IDF Vectorization
```

```
tfidf_vectorizer = TfidfVectorizer(max_features=5000)  # You can adjust the number of
features

X_train_tfidf = tfidf_vectorizer.fit_transform(X_train)

X_test_tfidf = tfidf_vectorizer.transform(X_test)


# Model Training - Naive Bayes Classifier

model = MultinomialNB()

model.fit(X_train_tfidf, y_train)


# Predict on the test set

y_pred = model.predict(X_test_tfidf)


# Model Evaluation

accuracy = accuracy_score(y_test, y_pred)

report = classification_report(y_test, y_pred)


print(f"Accuracy: {accuracy}")

print("Classification Report:")

print(report)
```

This example uses a basic TF-IDF vectorizer and a simple Multinomial Naive Bayes classifier. In practice, you can experiment with more sophisticated models and feature engineering techniques, and you should perform hyperparameter tuning to optimize the model's performance. Additionally, make sure to preprocess your text data, clean it, and handle any missing values as needed for your specific dataset.


•**Conclusion**

The passive-aggressive classifier performed the best here and gave an accuracy of 93.12%.


We can print a confusion matrix to gain insight into the number of false and true negatives and positives

Fake news detection techniques can be divided into those based on style and those based on content, or fact-checking. Too often it is assumed that bad style (bad spelling, bad punctuation, limited vocabulary, using terms of abuse, ungrammaticality, etc.) is a safe indicator of fake news.