



GOVERNMENT COLLEGE OF TECHNOLOGY

COIMBATORE

**(An Autonomous Institution Affiliated to Anna
University)**

COLLEGE CODE-7177

PROJECT TITLE: FAKE NEWS DETECTION USING NLP

TEAM MEMBERS:

SUBASH S -71772114146

ATCHAYA A -71772114104

ARUNTHATHI R -71772114103

RITHISHA P -71772114L01

PROBLEM STATEMENT:

It involves developing algorithms and models to accurately identify and classify misinformation in news articles or social media posts. The goal is to create a system that can distinguish between reliable and fake news sources based on various factors such as credibility, source reputation, fact-checking, and content analysis. The challenge lies in designing effective algorithms that can analyse large amounts of data and detect patterns or indicators of fake news with high accuracy.

PROBLEM SOLUTION :

- ❖ To solve the problem of fake news detection, various approaches can be taken.
- ❖ One solution involves using natural language processing (NLP) techniques to analyse the content of news articles and identify patterns that indicate misinformation.
- ❖ This can include analysing the language used, fact-checking claims, and identifying biased or misleading information.
- ❖ Additionally, machine learning algorithms can be trained on labelled datasets to classify news articles as either reliable or fake based on features such as source credibility, writing style, and social media engagement.
- ❖ Regular updates and improvements to the algorithms are crucial to keep up with evolving techniques used by those spreading fake news.
- ❖ To tackle the problem of fake news detection is by leveraging the power of crowdsourcing.
- ❖ Platforms can engage users to report and flag suspicious or misleading content, which can then be reviewed by fact-checkers or community moderators.
- ❖ This collaborative approach helps in identifying and verifying the accuracy of news articles through collective efforts.
- ❖ Additionally, educating users about media literacy and critical thinking skills can empower them to identify and question sources of fake news, thereby reducing its impact.

PROCESS

Step 1:

Choose the fake news dataset available on Kaggle, containing articles titles and text, along with their labels (genuine or fake).

Step 2: Data Preprocessing

- **Text Cleaning:** Remove special characters, numbers, and unnecessary whitespace from the text.
- **Tokenization:** Split the text into words or subwords.
- **Stopword Removal:** Remove common words like "and," "the," etc., which don't carry significant meaning.
- **Stemming/Lemmatization:** Reduce words to their root form for consistency.

Step 3: Feature Extraction

- **TF-IDF (Term Frequency-Inverse Document Frequency):** Convert text to numerical vectors, emphasizing words that are unique to specific documents.
- **Word Embeddings:** Utilize pre-trained word embeddings like Word2Vec, GloVe, or train your embeddings using algorithms like Word2Vec or FastText.
- **Sequence Padding:** Ensure all input sequences have the same length for processing.

Step 4: Model Selection

- **Recurrent Neural Networks (RNNs):** Especially LSTM (Long Short-Term Memory) networks can capture sequential patterns in text.
- **Bidirectional RNNs:** They learn from the sequence in both forward and backward directions, capturing context effectively.

- **Attention Mechanisms:** Enable the model to focus on different parts of the input sequence.
- **BERT (Bidirectional Encoder Representations from Transformers):** *For more advanced applications, you can use pre-trained transformer models like BERT for contextual embeddings.

Step 4: Model Training

- Split your data into training and testing sets.
- Train your selected model on the training data.
- Validate and fine-tune hyperparameters using techniques like cross-validation.
- Regularize your model to prevent overfitting (e.g., dropout layers).

Step 5: Model Evaluation

- Use appropriate metrics such as accuracy, precision, recall, and F1-score to evaluate your model's performance on the test set.
- Analyze false positives and false negatives to understand where the model fails.

Step 6: Iterate and Experiment

- Depending on the evaluation results, adjust your preprocessing steps, feature extraction methods, or even try different models.
- Experiment with ensemble methods if individual models aren't performing as desired.

Step 7: Deployment

- Once you have a well-performing model, deploy it in your desired environment. It could be a web application, API service, or any other platform where users can input text and get predictions.

CODE FOR EACH STEP

1. Load Data from Kaggle:

- Download the dataset from Kaggle and load it into a pandas DataFrame.

CODE:

```
import pandas as pd

# Load the dataset into a pandas DataFrame
df = pd.read_csv('kaggle_dataset.csv')
```

2. Text Cleaning and Preprocessing:

- Apply text cleaning and preprocessing functions to the 'text' column of the DataFrame.

CODE:

```
import re

from nltk.corpus import stopwords

from nltk.tokenize import word_tokenize
```

```
nltk.download('stopwords')
nltk.download('punkt')

def clean_text(text):
    text = re.sub(r"[^a-zA-Z]", " ", text) # Remove non-alphabetic
characters
    text = text.lower() # Convert text to lowercase
    words = word_tokenize(text) # Tokenize words
    words = [word for word in words if word not in
stopwords.words('english')] # Remove stopwords
    return ' '.join(words)

# Apply cleaning function to 'text' column
df['cleaned_text'] = df['text'].apply(clean_text)
```

3. TF-IDF Vectorization:

- Use scikit-learn's TF-IDF vectorizer to convert cleaned text into numerical vectors.

Code:

```
from sklearn.feature_extraction.text import TfidfVectorizer
```

```
# Initialize TF-IDF Vectorizer

vectorizer = TfidfVectorizer(max_features=5000) # Limit the
number of features


# Transform cleaned text into TF-IDF vectors

tfidf_matrix =
vectorizer.fit_transform(df['cleaned_text']).toarray()
```

4. Tokenization and Padding for Neural Networks:

- Use Keras for tokenization and padding sequences for neural network-based models.

CODE:

```
from tensorflow.keras.preprocessing.text import Tokenizer

from tensorflow.keras.preprocessing.sequence import
pad_sequences


# Tokenize and pad sequences

tokenizer = Tokenizer()

tokenizer.fit_on_texts(df['cleaned_text'])

sequences = tokenizer.texts_to_sequences(df['cleaned_text'])

padded_sequences = pad_sequences(sequences, maxlen=100,
padding='post')
```

Now, ``tfidf_matrix`` contains the TF-IDF vectors for your cleaned text, and ``padded_sequences`` contains the tokenized and padded sequences suitable for input into neural networks.

You can use these preprocessed data (``tfidf_matrix`` or ``padded_sequences``) to train machine learning models or neural networks for tasks like fake news detection. Remember to split your data into training and testing sets before training your models and choose appropriate algorithms/architectures based on your specific use case.

CONCLUSION:

After finishing all these process run this file and make api server with local host then paste the news . The compilation process will check whether the entered news is reliable or not.