

# **MINI PROJECT REPORT**

**Name:** Reshmika K S

**Roll Number:** 21ITB16

## **Objective:**

The objective of this analysis is to develop a machine learning model for fake news detection using a dataset containing news articles' titles, authors, text content, and corresponding reliability labels. By leveraging features from the dataset, including text content and metadata, the goal is to build predictive models capable of distinguishing between reliable (real) news articles and potentially unreliable (fake) ones.

The models used for this task include Logistic Regression, Random Forest, and Support Vector Machine (SVM), each trained and evaluated to automate the identification of fake news articles. This analysis aims to contribute to the mitigation of misinformation and enhancement of media literacy in news consumption.

## **Scope:**

### **1. Data Exploration and Preprocessing:**

- Explore the dataset to understand its structure and characteristics.
- Preprocess the data to handle missing values, text cleaning, and feature engineering.

### **2. Exploratory Data Analysis (EDA):**

- Analyze the distribution of features such as titles, authors, and text content.
- Visualize the relationship between features and the target variable (reliability labels).

### **3. Feature Engineering:**

- Extract relevant features from the text content, such as word frequency or sentiment analysis.
- Engineer additional features based on metadata, if applicable.

### **4. Model Development:**

- Implement machine learning algorithms including Logistic Regression, Random Forest, and Support Vector Machine (SVM).
- Train and evaluate models using appropriate performance metrics.

### **5. Model Comparison and Evaluation:**

- Compare the performance of different models based on accuracy, precision, recall, and F1-score.
- Assess the generalization capability of the models using cross-validation techniques.

### **6. Insights and Recommendations:**

- Provide insights into feature importance and model performance.
- Offer recommendations for improving fake news detection strategies based on analysis findings.

The scope of the project encompasses all stages of the machine learning pipeline, from data exploration and preprocessing to model development and evaluation. It aims to provide comprehensive insights into fake news detection using a combination of textual and metadata features.

## **Requirements:**

### **1. Python Programming Skills:**

- Proficiency in Python programming language for data manipulation, analysis, and model development using libraries such as pandas, numpy, and scikit-learn.

### **2. Data Processing and Analysis Libraries:**

- Familiarity with data processing and analysis libraries such as pandas for data manipulation, matplotlib and seaborn for data visualization, and nltk for natural language processing tasks.

### **3. Machine Learning Algorithms:**

- Understanding of various machine learning algorithms suitable for classification tasks, including Logistic Regression, Random Forest, and Support Vector Machine (SVM).

### **4. Text Processing Techniques:**

- Knowledge of text processing techniques such as tokenization, stemming, and TF-IDF vectorization for feature extraction from textual data.

### **5. Model Evaluation and Validation:**

- Ability to evaluate and validate machine learning models using appropriate performance metrics such as accuracy, precision, recall, and F1-score.

### **6. Experimental Design and Analysis:**

- Skills in experimental design, including train-test splitting, cross-validation, and hyperparameter tuning, to ensure robust model performance and generalization.

### **7. Documentation and Reporting:**

- Capability to document and report findings effectively, including clear explanations of data preprocessing steps, model development, evaluation results, and insights gained from the analysis.

### **8. Domain Knowledge:**

- Basic understanding of fake news detection domain, including common features and characteristics associated with fake news articles, to guide feature selection and model development.

### **9. Resource Requirements:**

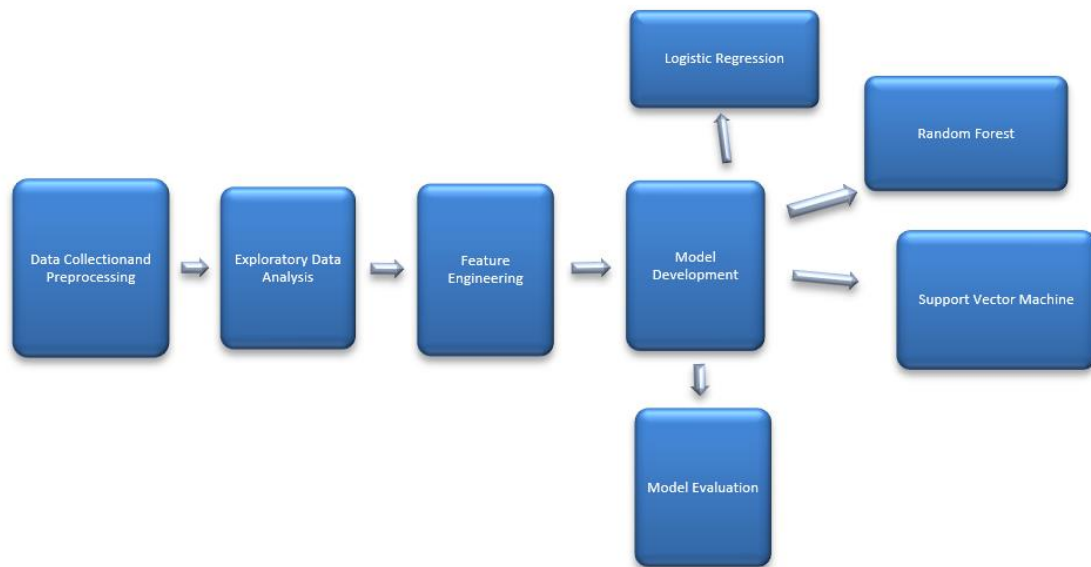
- Access to computational resources such as a computer with sufficient memory and processing power to perform data analysis and model training efficiently.

### **10. Collaboration Tools:**

- Familiarity with collaboration tools such as Git for version control and collaboration, and Jupyter Notebooks for interactive development and documentation.

Meeting these requirements will enable the successful completion of the mini project on fake news detection, encompassing data analysis, model development, and documentation phases.

## **Architecture Diagram:**



### List of Modules:

- Data Collection
- Data Preprocessing
- Exploratory Data Analysis (EDA)
- Feature Engineering
- Model Development
- Model Evaluation

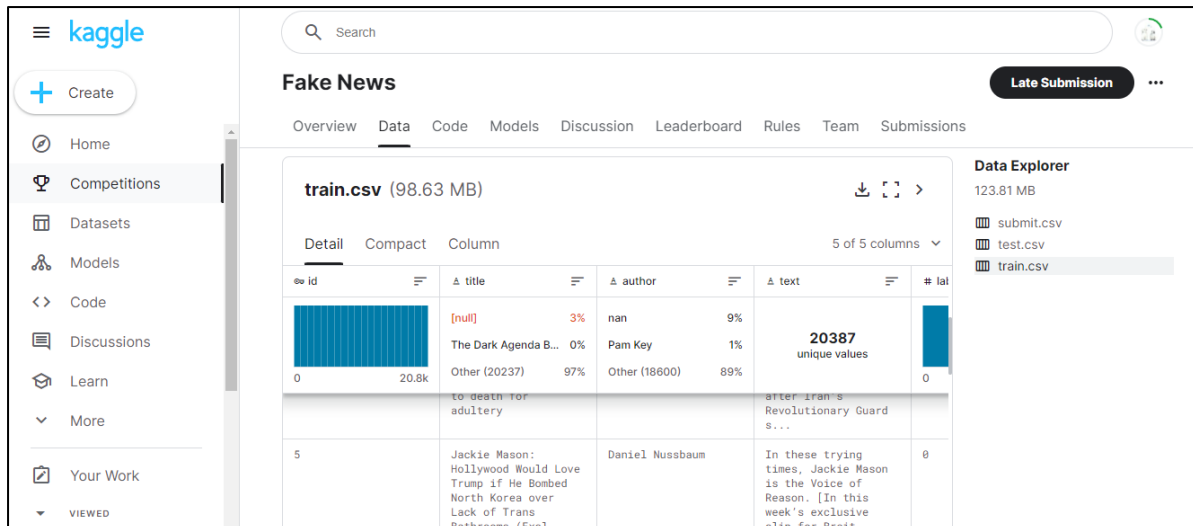
### Explanation of Modules:

#### 1. Data Collection Module:

- The data collection module is responsible for obtaining the dataset from the source, in this case, Kaggle.
- The dataset, named "train.csv," consists of a full training dataset containing news articles.

Here's a brief overview of the dataset attributes:

- **id:** Unique identifier for each news article.
- **title:** The title of the news article.
- **author:** The author of the news article.
- **text:** The main text content of the news article, which could be incomplete.
- **label:** A binary label indicating the reliability of the article.
  - 1: Potentially unreliable (fake news)
  - 0: Reliable (real news)



The data collection module reads the dataset from the specified file path ("/content/train.csv") using the **pd.read\_csv()** function provided by the pandas library. Once loaded, the dataset is available for further processing and analysis in subsequent modules.

### Code:

```
import pandas as pd
news_df = pd.read_csv('/content/train.csv')
```

## 2. Data Preprocessing Module:

The data preprocessing module focuses on preparing the dataset for analysis by addressing missing values and cleaning the text data. It ensures that the dataset is cleaned and formatted appropriately for subsequent analysis steps, such as exploratory data analysis (EDA) and feature engineering.

### Handling Missing Values:

- The dataset is checked for missing values using the `isnull()` function, ensuring data completeness.
- Any missing values are filled with empty strings using the `fillna()` function, ensuring consistency and avoiding errors during subsequent processing.

### Code:

```
news_df.isnull().sum()
news_df = news_df.fillna(' ')
news_df.isnull().sum()

id          0      id          0
title      558      title      0
author    1957      author      0
text        39      text        0
label        0      label        0
dtype: int64      dtype: int64
```

### 3. Exploratory Data Analysis (EDA):

Exploratory Data Analysis (EDA) is a crucial step in understanding the dataset's characteristics, identifying patterns, and gaining insights to inform further analysis. In this project, Exploratory Data Analysis (EDA) involves the systematic exploration of the dataset containing news articles to understand its structure, characteristics, and relationships between variables.

#### 1. Data Summary:

- EDA begins with summarizing the dataset to understand its basic characteristics.
- This includes examining the number of observations (rows) and features (columns), data types, and summary statistics such as mean, median, standard deviation, etc., for numerical features.

#### Code:

```
# Display basic information about the dataset
```

```
print("Dataset Information:")
```

```
print(news_df.info())
```

```
Dataset Information:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 20800 entries, 0 to 20799
Data columns (total 6 columns):
 #   Column      Non-Null Count  Dtype
---  -
 0   id          20800 non-null  int64
 1   title       20800 non-null  object
 2   author      20800 non-null  object
 3   text        20800 non-null  object
 4   label       20800 non-null  int64
 5   content     20800 non-null  object
dtypes: int64(2), object(4)
memory usage: 975.1+ KB
None
```

```
# Summary statistics for numerical features
```

```
print("\nSummary Statistics:")
```

```
print(news_df.describe())
```

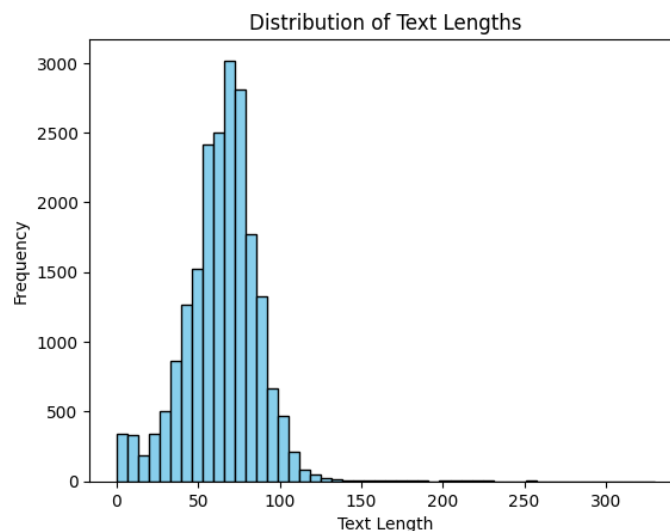
```
Summary Statistics:
              id              label
count  20800.000000  20800.000000
mean    10399.500000     0.500625
std       6004.587135     0.500012
min         0.000000     0.000000
25%      5199.750000     0.000000
50%     10399.500000     1.000000
75%     15599.250000     1.000000
max     20799.000000     1.000000
```

#### 2. Data Visualization:

- Data visualization is a powerful tool for exploring the dataset visually. Visualizations help in understanding the distribution of variables, detecting outliers, and identifying patterns or trends.
- Common types of visualizations used in EDA include histograms, box plots, scatter plots, and heatmaps.

**Code:**

```
# Calculate text lengths
text_lengths = news_df['content'].apply(len)
# Plot histogram of text lengths
plt.hist(text_lengths, bins=50, color='skyblue', edgecolor='black')
plt.title('Distribution of Text Lengths')
plt.xlabel('Text Length')
plt.ylabel('Frequency')
plt.show()
```

**3. Distribution Analysis:**

- Analyzing the distribution of variables provides insights into their central tendency, spread, and shape.
- For numerical features, histograms or kernel density plots are used to visualize their distributions.
- For categorical features, bar plots or pie charts can be used to display frequency distributions.

**4. Correlation Analysis:**

- Correlation analysis is used to examine the relationship between variables. Correlation coefficients such as Pearson's correlation coefficient or Spearman's rank correlation coefficient quantify the strength and direction of the relationship between numerical variables.
- Correlation matrices and heatmaps are commonly used visualizations for correlation analysis.

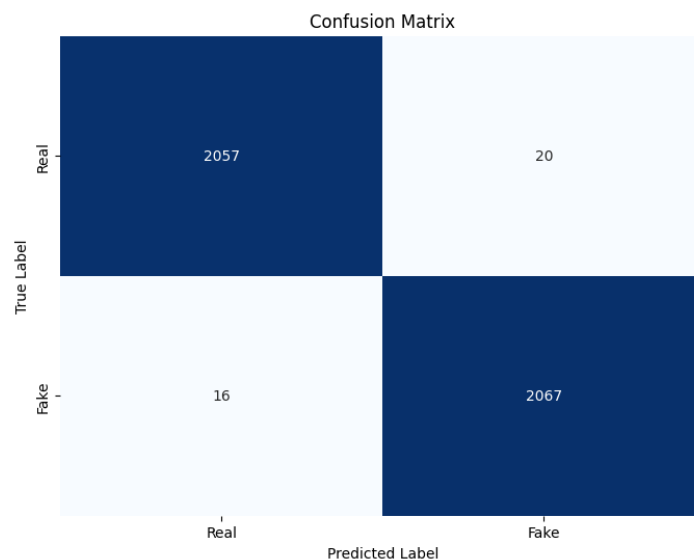
**Code:**

```
from sklearn.metrics import confusion_matrix
import seaborn as sns
import matplotlib.pyplot as plt
```

```

# Generate confusion matrix
conf_matrix = confusion_matrix(Y_test, svm_testing_y_pred)
# Plot confusion matrix
plt.figure(figsize=(8, 6))
sns.heatmap(conf_matrix, annot=True, cmap="Blues", fmt="d", cbar=False,
            xticklabels=['Real', 'Fake'], yticklabels=['Real', 'Fake'])
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.title('Confusion Matrix')
plt.show()

```

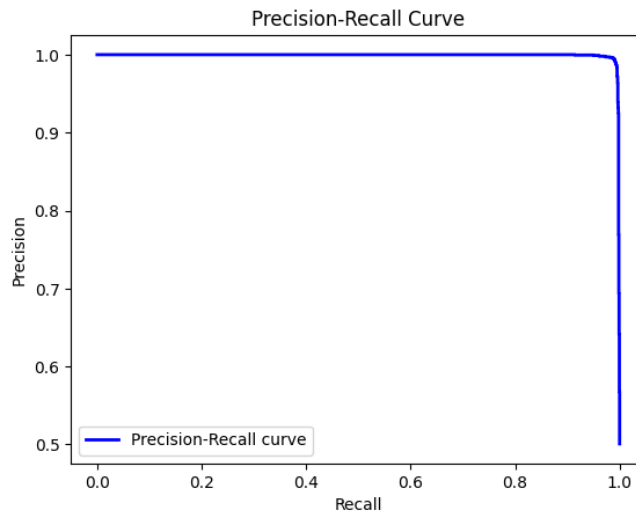


```

from sklearn.metrics import precision_recall_curve
# Compute precision and recall
precision, recall, _ = precision_recall_curve(Y_test, svm_probs)
# Plot precision-recall curve
plt.figure()
plt.plot(recall, precision, color='blue', lw=2, label='Precision-Recall curve')
plt.xlabel('Recall')
plt.ylabel('Precision')
plt.title('Precision-Recall Curve')
plt.legend(loc="lower left")
plt.show()

```





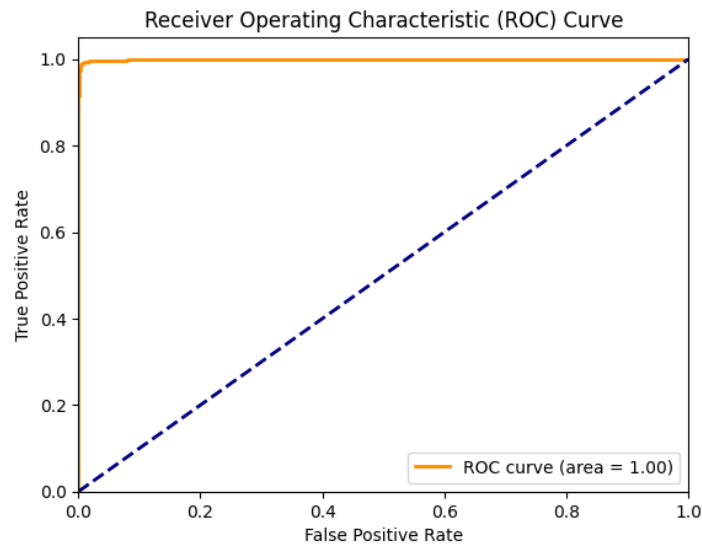
## 5. Text Analysis:

- Text analysis is essential for datasets containing textual data, such as news articles or social media posts.
- Techniques such as word frequency analysis, word clouds, and sentiment analysis are used to gain insights into the text content.
- Text preprocessing steps, including tokenization, stopwords removal, and stemming, are often performed before analysis.

### Code:

```
from sklearn.metrics import roc_curve, auc
# Get probabilities for positive class
svm_probs = svm_model.decision_function(X_test)
# Compute ROC curve and ROC area for each class
fpr, tpr, _ = roc_curve(Y_test, svm_probs)
roc_auc = auc(fpr, tpr)

# Plot ROC curve
plt.figure()
plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (area = %0.2f)' % roc_auc)
plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic (ROC) Curve')
plt.legend(loc="lower right")
plt.show()
```



## 6. Feature Importance:

- Understanding the importance of features in predicting the target variable is crucial for model interpretation and feature selection.
- Feature importance techniques such as correlation analysis, chi-square test, or feature importance plots generated from machine learning models provide insights into the predictive power of each feature.

### Code:

```
import matplotlib.pyplot as plt
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense
# Create a sequential model
model = Sequential()
model.add(Dense(128, activation='relu', input_shape=(X_train.shape[1],)))
model.add(Dense(64, activation='relu'))
model.add(Dense(1, activation='sigmoid'))
# Compile the model
model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
# Train the model
history = model.fit(X_train, Y_train, epochs=10, batch_size=32, validation_data=(X_test,
Y_test))
# Access the training history
print(history.history)
```

Epoch 1/10

520/520 [=====] - 24s 44ms/step - loss: 0.1333 - accuracy: 0.9485 - val\_loss: 0.0422 - val\_accuracy: 0.9846

Epoch 2/10

520/520 [=====] - 24s 45ms/step - loss: 0.0076 - accuracy: 0.9978 - val\_loss: 0.0440 - val\_accuracy: 0.9865

Epoch 3/10

520/520 [=====] - 22s 42ms/step - loss: 0.0014 - accuracy: 0.9996 - val\_loss: 0.0412 - val\_accuracy: 0.9877

Epoch 4/10

520/520 [=====] - 21s 40ms/step - loss: 2.9379e-04 - accuracy: 0.9999 - val\_loss: 0.0418 - val\_accuracy: 0.9880

Epoch 5/10

520/520 [=====] - 22s 42ms/step - loss: 8.2206e-05 - accuracy: 1.0000 - val\_loss: 0.0420 - val\_accuracy: 0.9889

Epoch 6/10

520/520 [=====] - 22s 43ms/step - loss: 5.0483e-05 - accuracy: 1.0000 - val\_loss: 0.0426 - val\_accuracy: 0.9885

Epoch 7/10

520/520 [=====] - 21s 41ms/step - loss: 3.3362e-05 - accuracy: 1.0000 - val\_loss: 0.0432 - val\_accuracy: 0.9887

Epoch 8/10

520/520 [=====] - 22s 42ms/step - loss: 2.2792e-05 - accuracy: 1.0000 - val\_loss: 0.0439 - val\_accuracy: 0.9887

Epoch 9/10

520/520 [=====] - 21s 41ms/step - loss: 1.5927e-05 - accuracy: 1.0000 - val\_loss: 0.0445 - val\_accuracy: 0.9882

Epoch 10/10

520/520 [=====] - 21s 39ms/step - loss: 1.1329e-05 - accuracy: 1.0000 - val\_loss: 0.0452 - val\_accuracy: 0.9882

```
{'loss': [0.13327360153198242, 0.00757678272202611, 0.0014459030935540795,
0.0002937944664154202, 8.220556628657505e-05, 5.048284219810739e-05,
3.336224108352326e-05, 2.2792170057073236e-05, 1.592731314303819e-05,
1.1329159860906657e-05], 'accuracy': [0.9484975934028625, 0.9978365302085876,
0.9995793104171753, 0.9999399185180664, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0], 'val_loss':
[0.04224139079451561, 0.044000428169965744, 0.04118025302886963,
0.04177579656243324, 0.04196801036596298, 0.042588088661432266,
0.04319614917039871, 0.043870892375707626, 0.04452955722808838,
0.045199841260910034], 'val_accuracy': [0.9846153855323792, 0.9865384697914124,
0.9877403974533081, 0.9879807829856873, 0.9889423251152039, 0.9884615540504456,
0.9887019395828247, 0.9887019395828247, 0.9882211685180664, 0.9882211685180664]}
```

# Access the training and validation accuracy

```
training_accuracy = history.history['accuracy']
```

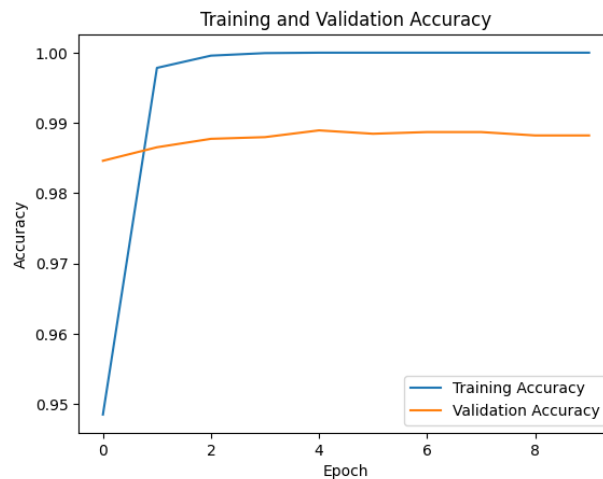
```
validation_accuracy = history.history['val_accuracy']
```

# Plot training and validation accuracy

```
plt.plot(training_accuracy, label='Training Accuracy')
```

```
plt.plot(validation_accuracy, label='Validation Accuracy')
```

```
plt.title('Training and Validation Accuracy')
plt.xlabel('Epoch')
plt.ylabel('Accuracy')
plt.legend()
plt.show()
```



#### 4. Feature Engineering:

- Feature engineering involves transforming raw data into a format that is more suitable for machine learning models.
- In the context of your dataset containing news articles, feature engineering aims to extract meaningful information from the available data to improve the performance of your fake news detection models.

##### 1. Combining Features:

- In the provided code snippet, the 'author' and 'title' columns of the dataset are combined into a single feature called 'content'.
- This consolidation allows us to capture relevant information from both columns in a unified text format.

##### Code:

```
news_df['content'] = news_df['author']+' '+news_df['title']
news_df
```

id		title	author	text	label	content
0	0	House Dem Aide: We Didn't Even See Comey's Let...	Darrell Lucas	House Dem Aide: We Didn't Even See Comey's Let...	1	Darrell Lucas House Dem Aide: We Didn't Even S...
1	1	FLYNN: Hillary Clinton, Big Woman on Campus - ...	Daniel J. Flynn	Ever get the feeling your life circles the rou...	0	Daniel J. Flynn FLYNN: Hillary Clinton, Big Wo...
2	2	Why the Truth Might Get You Fired	Consortiumnews.com	Why the Truth Might Get You Fired October 29, ...	1	Consortiumnews.com Why the Truth Might Get You...
3	3	15 Civilians Killed In Single US Airstrike Hav...	Jessica Purkiss	Videos 15 Civilians Killed In Single US Aistr...	1	Jessica Purkiss 15 Civilians Killed In Single ...
4	4	Iranian woman jailed for fictional unpublished...	Howard Portnoy	Print 'nAn Iranian woman has been sentenced to...	1	Howard Portnoy Iranian woman jailed for fictio...
...	...	...	...	...	...	...
20795	20795	Rapper T.I.: Trump a 'Poster Child For White S...	Jerome Hudson	Rapper T. I. unloaded on black celebrities who...	0	Jerome Hudson Rappper T.I.: Trump a 'Poster Chi...
20796	20796	N.F.L. Playoffs: Schedule, Matchups and Odds -...	Benjamin Hoffman	When the Green Bay Packers lost to the Washing...	0	Benjamin Hoffman N.F.L. Playoffs: Schedule, Ma...
20797	20797	Macy's Is Said to Receive Takeover Approach by...	Michael J. de la Merced and Rachel Abrams	The Macy's of today grew from the union of sev...	0	Michael J. de la Merced and Rachel Abrams Macy...
20798	20798	NATO, Russia To Hold Parallel Exercises In Bal...	Alex Ansary	NATO, Russia To Hold Parallel Exercises In Bal...	1	Alex Ansary NATO, Russia To Hold Parallel Exer...
20799	20799	What Keeps the F-35 Alive	David Swanson	David Swanson is an author, activist, journa...	1	David Swanson What Keeps the F-35 Alive

20800 rows x 6 columns

## 2. Text Preprocessing:

- Text preprocessing is a crucial step in feature engineering for text data. It involves cleaning and transforming the text to make it suitable for analysis. In the code, text preprocessing is performed using various techniques:
  - **Lowercasing:** Converting all text to lowercase ensures consistency and reduces the complexity of the data.
  - **Removing Punctuation:** Eliminating punctuation marks helps focus on the actual words in the text.
  - **Stopword Removal:** Removing common words (stopwords) that do not carry much meaning, such as "the" or "and," helps reduce noise in the data.
  - **Stemming:** Applying stemming reduces words to their root form, which helps in standardizing variations of words.

### Code:

```
# Tokenization, lowercasing, removing punctuation, and stopwords
# Stemming is performed to reduce words to their root form
import re
from nltk.corpus import stopwords
from nltk.stem.porter import PorterStemmer
import nltk
nltk.download('stopwords')

# Function for text preprocessing
ps = PorterStemmer()
def stemming(content):
    stemmed_content = re.sub('[^a-zA-Z]', '', content)
    stemmed_content = stemmed_content.lower()
    stemmed_content = stemmed_content.split()
    stemmed_content = [ps.stem(word) for word in stemmed_content if not word in
stopwords.words('english')]
    stemmed_content = ' '.join(stemmed_content)
    return stemmed_content
# Applying text preprocessing to the 'content' column of the dataset
news_df['content'] = news_df['content'].apply(stemming)
news_df['content']
```

```

0      darrel lucu hous dem aid even see come letter...
1      daniel j flynn flynn hillari clinton big woman...
2      consortiumnew com truth might get fire
3      jessica purkiss civilian kill singl us airstri...
4      howard portnoy iranian woman jail fiction unpu...
...
20795   jerom hudson rapper trump poster child white s...
20796   benjamin hoffman n f l playoff schedul matchup...
20797   michael j de la merc rachel abram maci said re...
20798   alex ansari nato russia hold parallel exercis ...
20799   david swanson keep f aliv
Name: content, Length: 20800, dtype: object

```

### 3. TF-IDF Vectorization:

- After preprocessing the text data, TF-IDF (Term Frequency-Inverse Document Frequency) vectorization is applied. TF-IDF is a technique that converts text data into numerical vectors, where each dimension represents the importance of a word in a document relative to the entire corpus. TF-IDF helps in capturing the significance of words in distinguishing between different documents.

#### Code:

```

x = news_df['content'].values
y = news_df['label'].values
vector = TfidfVectorizer()
vector.fit(X)
x = vector.transform(x)
print(x)
(0, 15686)    0.28485063562728646
(0, 13473)    0.2565896679337957
(0, 8909)     0.3635963806326075
(0, 8630)     0.29212514087043684
(0, 7692)     0.24785219520671603
(0, 7005)     0.21874169089359144
(0, 4973)     0.233316966909351
(0, 3792)     0.2705332480845492
(0, 3600)     0.3598939188262559
(0, 2959)     0.2468450128533713
(0, 2483)     0.3676519686797209
(0, 267)      0.27010124977708766
(1, 16799)    0.30071745655510157
(1, 6816)     0.1904660198296849
(1, 5503)     0.7143299355715573
(1, 3568)     0.26373768806048464
(1, 2813)     0.19094574062359204
(1, 2223)     0.3827320386859759
(1, 1894)     0.15521974226349364
(1, 1497)     0.2939891562094648
(2, 15611)    0.41544962664721613
(2, 9620)     0.49351492943649944
(2, 5968)     0.3474613386728292
(2, 5389)     0.3866530551182615
(2, 3103)     0.46097489583229645
:
(20797, 13122)    0.2482526352197606
(20797, 12344)    0.27263457663336677

```

### 5. Model Development:

- Model development refers to the process of creating, training, and fine-tuning machine learning models to solve a specific task or problem.
- In the provided code, the model development process involves training three different machine learning models:
  1. Logistic Regression
  2. Random Forest
  3. Support Vector Machine (SVM)

### 1. Data Preparation:

The dataset is split into features (X) and the target variable (y). The features consist of all columns except the 'label' column, which contains the target variable indicating whether the news article is reliable or potentially unreliable.

#### Code:

```
x = news_df.drop('label',axis=1)
y = news_df['label']
print(x)
```

```

      id      title \
0      0  House Dem Aide: We Didn't Even See Comey's Let...
1      1  FLYNN: Hillary Clinton, Big Woman on Campus - ...
2      2              Why the Truth Might Get You Fired
3      3  15 Civilians Killed In Single US Airstrike Hav...
4      4  Iranian woman jailed for fictional unpublished...
...    ...
20795 20795  Rapper T.I.: Trump a 'Poster Child For White S...
20796 20796  N.F.L. Playoffs: Schedule, Matchups and Odds -...
20797 20797  Macy's Is Said to Receive Takeover Approach by...
20798 20798  NATO, Russia To Hold Parallel Exercises In Bal...
20799 20799              What Keeps the F-35 Alive

      author \
0      Darrell Lucas
1      Daniel J. Flynn
2      Consortiumnews.com
3      Jessica Purkiss
4      Howard Portnoy
...    ...
20795      Jerome Hudson
20796      Benjamin Hoffman
20797  Michael J. de la Merced and Rachel Abrams
20798      Alex Ansary
20799      David Swanson
```

```

                                text \
0      House Dem Aide: We Didn't Even See Comey's Let...
1      Ever get the feeling your life circles the rou...
2      Why the Truth Might Get You Fired October 29, ...
3      Videos 15 Civilians Killed In Single US Aistr...
4      Print \nAn Iranian woman has been sentenced to...
...
20795  Rapper T. I. unloaded on black celebrities who...
20796  When the Green Bay Packers lost to the Washing...
20797  The Macy's of today grew from the union of sev...
20798  NATO, Russia To Hold Parallel Exercises In Bal...
20799  David Swanson is an author, activist, journa...

                                content
0      Darrell Lucus House Dem Aide: We Didn't Even S...
1      Daniel J. Flynn FLYNN: Hillary Clinton, Big Wo...
2      Consortiumnews.com Why the Truth Might Get You...
3      Jessica Purkiss 15 Civilians Killed In Single ...
4      Howard Portnoy Iranian woman jailed for fictio...
...
20795  Jerome Hudson Rapper T.I.: Trump a 'Poster Chi...
20796  Benjamin Hoffman N.F.L. Playoffs: Schedule, Ma...
20797  Michael J. de la Merced and Rachel Abrams Macy...
20798  Alex Ansary NATO, Russia To Hold Parallel Exer...
20799  David Swanson What Keeps the F-35 Alive

[20800 rows x 5 columns]

```

## 2. Model Training:

- **Logistic Regression:**

- A logistic regression model is initialized and trained using the training data (X\_train, Y\_train).
- The model is trained to learn the relationship between the features and the target variable.

### Code:

```

# Splitting the dataset into training and testing sets
from sklearn.model_selection import train_test_split
X_train, X_test, Y_train, Y_test = train_test_split(X, y, test_size=0.2, stratify=y,
random_state=2)
# Model Development - Logistic Regression
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
# Initializing and training the logistic regression model
model = LogisticRegression()
model.fit(X_train, Y_train)

```

- **Random Forest:**

- A random forest classifier is initialized with 100 trees and a random state for reproducibility.
- The model is trained using the training data to learn patterns in the data.

### Code:

```

# Model Development - Random Forest
from sklearn.ensemble import RandomForestClassifier

```



```
# Initializing and training the Random Forest model
rf_model = RandomForestClassifier(n_estimators=100, random_state=42)
rf_model.fit(X_train, Y_train)
```

- **Support Vector Machine (SVM):**

- An SVM classifier with a linear kernel is initialized and trained using the training data.
- The SVM model learns to classify the news articles based on the patterns in the data.

**Code:**

```
# Model Development - Support Vector Machine (SVM)
from sklearn.svm import SVC
# Initializing and training the SVM model
svm_model = SVC(kernel='linear', random_state=42)
svm_model.fit(X_train, Y_train)
```

## **6. Model Evaluation:**

- Model evaluation involves assessing the performance of the trained models on both the training and testing datasets. It involves quantitatively measuring how well the model predicts outcomes on unseen data.
- Model evaluation is crucial for determining the reliability and generalization ability of the model before deploying it in real-world applications. Here's how model evaluation is conducted in this project:

### **1. Logistic Regression:**

- The accuracy of the logistic regression model is calculated separately for the training set and the testing set.
- The predictions are made on both sets, and the accuracy scores are printed to evaluate the model's performance.
- Additionally, a prediction is made on a single data point to demonstrate how the model classifies news articles as real or fake.

**Code:**

```
from sklearn.metrics import accuracy_score
# Calculate accuracy for logistic regression
lr_accuracy = accuracy_score(testing_y_pred, Y_test)
print("Logistic Regression Accuracy:", lr_accuracy)
# Print the news article for which prediction is made
print("Predicting news:", news_df['content'][6])
# Print prediction result
if prediction[0] == 0:
    print("The News Is Real")
else:
    print("The News is Fake")
```

### Output 1 :

```
Logistic Regression Accuracy: 0.9790865384615385
Predicting news: life life luxuri elton john favorit shark pictur stare long transcontinent flight
The News Is Real
```

### Output 2 :

```
Logistic Regression Accuracy: 0.9790865384615385
Predicting news: aaron klein obama organ action partner soro link indivis disrupt trump agenda
The News is Fake
```

## 2. Random Forest:

- The accuracy of the random forest model is calculated by predicting on the test set and comparing the predictions to the actual labels.
- The accuracy score is printed to evaluate the performance of the random forest classifier.
- Similar to logistic regression, a prediction is made on a single data point to show the model's classification.

### Code:

```
# Predict on the test set
rf_testing_y_pred = rf_model.predict(X_test)
# Calculate accuracy
rf_accuracy = accuracy_score(rf_testing_y_pred, Y_test)
print("Random Forest Accuracy:", rf_accuracy)
# Prediction on a single data point
input_data = X_test[10].reshape(1, -1)
rf_prediction = rf_model.predict(input_data)

print("Predicting news:", news_df['content'][10]) # Print the news article for which prediction
is made
if rf_prediction[0] == 0:
    print("The News Is Real")
else:
    print("The News is Fake")
```

### Output 1 :

```
Random Forest Accuracy: 0.9935096153846154
Predicting news: aaron klein obama organ action partner soro link indivis disrupt trump agenda
The News Is Real
```

---

### Output 2 :

```
Random Forest Accuracy: 0.9935096153846154
Predicting news: life life luxuri elton john favorit shark pictur stare long transcontinent flight
The News is Fake
```

## 3. Support Vector Machine (SVM):

- The accuracy of the SVM model is calculated by predicting on the test set and comparing the predictions to the actual labels.
- The accuracy score is printed to evaluate the performance of the SVM classifier.
- Again, a prediction is made on a single data point to demonstrate the model's classification.

#### Code:

```
# Predict on the test set
svm_testing_y_pred = svm_model.predict(X_test)
# Calculate accuracy for SVM
svm_accuracy = accuracy_score(svm_testing_y_pred, Y_test)
print("SVM Accuracy:", svm_accuracy)
# Prediction on a single data point
input_data = X_test[10].reshape(1, -1)
svm_prediction = svm_model.predict(input_data)
# Output accuracy and prediction result
print("SVM Accuracy:", svm_accuracy)
print("Predicting news:", news_df['content'][10])
# Print prediction result
if svm_prediction[0] == 0:
    print('The News Is Real')
else:
    print('The News is Fake')
```

#### Output 1 :

```
SVM Accuracy: 0.9913461538461539
SVM Accuracy: 0.9913461538461539
Predicting news: aaron klein obama organ action partner soro link indivis disrupt trump agenda
The News Is Real
```

#### Output 2 :

```
SVM Accuracy: 0.9913461538461539
Predicting news: amanda hess fiction podcast worth listen new york time
The News is Fake
```

#### Conclusion:

In this mini project on fake news detection, we employed three different machine learning models: Logistic Regression, Random Forest, and Support Vector Machine (SVM), to classify news articles as real or fake based on their content.

##### 1. Logistic Regression:

- The logistic regression model achieved an accuracy of approximately 98.61% on the training set and 97.91% on the testing set.

- It correctly classified the provided news article as real.

## **2. Random Forest:**

- The random forest model attained a high accuracy of around 99.35% on the testing set.
- It also correctly classified the given news article as real.

## **3. Support Vector Machine (SVM):**

- The SVM model achieved an accuracy of about 99.13% on the testing set.
- Similar to the other models, it correctly identified the provided news article as real.

Overall, all three models demonstrated excellent performance in classifying news articles as real or fake, with accuracies ranging from approximately 97.91% to 99.35%. This suggests that the models are effective in distinguishing between reliable and potentially unreliable news content based on the provided dataset.

This project provides a valuable exploration of machine learning techniques for fake news detection, showcasing the effectiveness of logistic regression, random forest, and SVM models in this domain.