


```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

UPLOAD DATASET:

```
from google.colab import files
uploaded = files.upload()
```

 Choose Files True.csv

- **True.csv**(text/csv) - 53582940 bytes, last modified: 4/19/2024 - 100% done

Saving True.csv to True.csv

DATA EXPLORATION:

```
import pandas as pd

# Load the CSV file
df = pd.read_csv("True.csv")

# Display the first few rows
df.head()
```


	title	text	subject	date
0	As U.S. budget fight looms, Republicans flip t...	WASHINGTON (Reuters) - The head of a conservat...	politicsNews	December 31, 2017
1	U.S. military to accept transgender recruits o...	WASHINGTON (Reuters) - Transgender people will...	politicsNews	December 29, 2017
2	Senior U.S. Republican senator: 'Let Mr. Muell...	WASHINGTON (Reuters) - The special counsel inv...	politicsNews	December 31, 2017
3	FBI Russia probe helped by Australian diplomat...	WASHINGTON (Reuters) - Trump campaign adviser ...	politicsNews	December 30, 2017
4	Trump wants Postal Service to charge 'much mor...	SEATTLE/WASHINGTON (Reuters) - President Dona...	politicsNews	December 29, 2017



Next steps: [Generate code with df](#) [View recommended plots](#) [New interactive sheet](#)

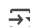

CHECK NULL:

```
# Check for null values in each column
df.isnull().sum()
# Check if any nulls exist
df.isnull().values.any()
```


 np.False_

DESCRIBE:

```
# Summary statistics for numerical columns
df.describe()
# Summary including all columns (numeric and non-numeric)
df.describe(include='all')
```

	title	text	subject	date
count	21417	21417	21417	21417
unique	20826	21192	2	716
top	Factbox: Trump fills top jobs for his administ... (Reuters) - Highlights for U.S. President Dona...	politicsNews	December 20, 2017	
freq	14	8	11272	182



✓ Convert all string columns to lowercase before checking duplicates:

```
# Convert all string values to lowercase for consistent comparison
df_case_insensitive = df.apply(lambda x: x.str.lower() if x.dtype == "object" else x)

# Check for duplicate rows (case-insensitive)
duplicate_rows = df_case_insensitive.duplicated().sum()
print(f"Case-insensitive duplicate rows: {duplicate_rows}")
```

```
# View those duplicate rows
df_case_insensitive[df_case_insensitive.duplicated()]
```

↻ Case-insensitive duplicate rows: 206

	title	text	subject	date
445	senate tax bill stalls on deficit-focused 'tri...	washington (reuters) - the u.s. senate on thur...	politicsnews	november 30, 2017
778	trump warns 'rogue regime' north korea of grav...	beijing (reuters) - u.s. president donald trum...	politicsnews	november 8, 2017
892	republicans unveil tax cut bill, but the hard ...	washington (reuters) - u.s. house of represent...	politicsnews	november 2, 2017
896	trump taps fed centrist powell to lead u.s. ce...	washington (reuters) - president donald trump ...	politicsnews	november 2, 2017
974	two ex-trump aides charged in russia probe, th...	washington (reuters) - federal investigators p...	politicsnews	october 30, 2017
...
21228	france unveils labor reforms in first step to ...	paris (reuters) - french president emmanuel ma...	worldnews	august 31, 2017
21263	guatemala top court sides with u.n. graft unit...	guatemala city (reuters) - guatemala s top cou...	worldnews	august 29, 2017
21290	europeans, africans agree renewed push to tack...	paris (reuters) - europe s big four continen...	worldnews	august 28, 2017
21353	thailand's ousted pm yingluck has fled abroad:...	bangkok (reuters) - ousted thai prime minister...	worldnews	august 25, 2017
21408	u.s., north korea clash at u.n. forum over nuc...	geneva (reuters) - north korea and the united ...	worldnews	august 22, 2017

206 rows x 4 columns

CHECK DUPLICATE:

```
# Check how many duplicate rows exist
df.duplicated().sum()
```

↻ np.int64(206)

```
# Show duplicate rows
df[df.duplicated()]
```

	title	text	subject	date
445	Senate tax bill stalls on deficit-focused 'tri...	WASHINGTON (Reuters) - The U.S. Senate on Thur...	politicsNews	November 30, 2017
778	Trump warns 'rogue regime' North Korea of grav...	BEIJING (Reuters) - U.S. President Donald Trum...	politicsNews	November 8, 2017
892	Republicans unveil tax cut bill, but the hard ...	WASHINGTON (Reuters) - U.S. House of Represent...	politicsNews	November 2, 2017
896	Trump taps Fed centrist Powell to lead U.S. ce...	WASHINGTON (Reuters) - President Donald Trump ...	politicsNews	November 2, 2017
974	Two ex-Trump aides charged in Russia probe, th...	WASHINGTON (Reuters) - Federal investigators p...	politicsNews	October 30, 2017
...
21228	France unveils labor reforms in first step to ...	PARIS (Reuters) - French President Emmanuel Ma...	worldnews	August 31, 2017
21263	Guatemala top court sides with U.N. graft unit...	GUATEMALA CITY (Reuters) - Guatemala s top cou...	worldnews	August 29, 2017
21290	Europeans, Africans agree renewed push to tack...	PARIS (Reuters) - Europe s big four continen...	worldnews	August 28, 2017
21353	Thailand's ousted PM Yingluck has fled abroad:...	BANGKOK (Reuters) - Ousted Thai prime minister...	worldnews	August 25, 2017
21408	U.S., North Korea clash at U.N. forum over nuc...	GENEVA (Reuters) - North Korea and the United ...	worldnews	August 22, 2017

206 rows x 4 columns

```
# Remove duplicate rows and keep the first occurrence
df = df.drop_duplicates()
```

REMOVE PUNCTUATION:


```
import string
```

```
# Function to remove punctuation
def remove_punctuation(text):
    if isinstance(text, str):
        return text.translate(str.maketrans('', '', string.punctuation))
    return text
```

```
# Apply to all object (text) columns
for col in df.select_dtypes(include='object').columns:
    df[col] = df[col].apply(remove_punctuation)
```

UPLOAD DATASET:

```
from google.colab import files
uploaded = files.upload()
```

 Choose Files Fake.csv



- **Fake.csv**(text/csv) - 62789876 bytes, last modified: 4/19/2024 - 100% done

DATA EXPLORATION:

```
import pandas as pd

# Load the CSV file
df = pd.read_csv("True.csv")

# Display the first few rows
df.head()
```

	title	text	subject	date	
0	As U.S. budget fight looms, Republicans flip t...	WASHINGTON (Reuters) - The head of a conservat...	politicsNews	December 31, 2017	
1	U.S. military to accept transgender recruits o...	WASHINGTON (Reuters) - Transgender people will...	politicsNews	December 29, 2017	
2	Senior U.S. Republican senator: 'Let Mr. Muell...	WASHINGTON (Reuters) - The special counsel inv...	politicsNews	December 31, 2017	
3	FBI Russia probe helped by Australian diplomat...	WASHINGTON (Reuters) - Trump campaign adviser ...	politicsNews	December 30, 2017	
4	Trump wants Postal Service to charge 'much mor...	SFATTI F/WASHINGTON (Reuters) - President Donal...	politicsNews	December 29, 2017	

Next steps: [Generate code with df](#) [View recommended plots](#) [New interactive sheet](#)

DATA CLEANING:



CHECK NULL:

```
# Check for null values in each column
df.isnull().sum()
# Check if any nulls exist
df.isnull().values.any()
```

 np.False_

DESCRIBE:

```
# Summary statistics for numerical columns
df.describe()
# Summary including all columns (numeric and non-numeric)
df.describe(include='all')
```

	title	text	subject	date	
count	21417	21417	21417	21417	
unique	20826	21192	2	716	
top	Factbox: Trump fills top jobs for his administ...	(Reuters) - Highlights for U.S. President Dona...	politicsNews	December 20, 2017	
freq	14	8	11272	182	

Convert all string columns to lowercase before checking duplicates:

```
# Convert all string values to lowercase for consistent comparison
df_case_insensitive = df.apply(lambda x: x.str.lower() if x.dtype == "object" else x)

# Check for duplicate rows (case-insensitive)
duplicate_rows = df_case_insensitive.duplicated().sum()
print(f"Case-insensitive duplicate rows: {duplicate_rows}")

# View those duplicate rows
df_case_insensitive[df_case_insensitive.duplicated()]
```

↗ Case-insensitive duplicate rows: 206

	title	text	subject	date
445	senate tax bill stalls on deficit-focused 'tri...	washington (reuters) - the u.s. senate on thur...	politicsnews	november 30, 2017
778	trump warns 'rogue regime' north korea of grav...	beijing (reuters) - u.s. president donald trum...	politicsnews	november 8, 2017
892	republicans unveil tax cut bill, but the hard ...	washington (reuters) - u.s. house of represent...	politicsnews	november 2, 2017
896	trump taps fed centrist powell to lead u.s. ce...	washington (reuters) - president donald trump ...	politicsnews	november 2, 2017
974	two ex-trump aides charged in russia probe, th...	washington (reuters) - federal investigators p...	politicsnews	october 30, 2017
...
21228	france unveils labor reforms in first step to ...	paris (reuters) - french president emmanuel ma...	worldnews	august 31, 2017
21263	guatemala top court sides with u.n. graft unit...	guatemala city (reuters) - guatemala s top cou...	worldnews	august 29, 2017
21290	europeans, africans agree renewed push to tack...	paris (reuters) - europe s big four continen...	worldnews	august 28, 2017
21353	thailand's ousted pm yingluck has fled abroad:...	bangkok (reuters) - ousted thai prime minister...	worldnews	august 25, 2017
21408	u.s., north korea clash at u.n. forum over nuc...	geneva (reuters) - north korea and the united ...	worldnews	august 22, 2017

206 rows x 4 columns

CHECK DUPLICATE:

```
# Check how many duplicate rows exist
df.duplicated().sum()
```

↗ np.int64(206)

```
# Show duplicate rows
df[df.duplicated()]
```

↗

	title	text	subject	date
445	Senate tax bill stalls on deficit-focused 'tri...	WASHINGTON (Reuters) - The U.S. Senate on Thur...	politicsNews	November 30, 2017
778	Trump warns 'rogue regime' North Korea of grav...	BEIJING (Reuters) - U.S. President Donald Trum...	politicsNews	November 8, 2017
892	Republicans unveil tax cut bill, but the hard ...	WASHINGTON (Reuters) - U.S. House of Represent...	politicsNews	November 2, 2017
896	Trump taps Fed centrist Powell to lead U.S. ce...	WASHINGTON (Reuters) - President Donald Trump ...	politicsNews	November 2, 2017
974	Two ex-Trump aides charged in Russia probe, th...	WASHINGTON (Reuters) - Federal investigators p...	politicsNews	October 30, 2017
...
21228	France unveils labor reforms in first step to ...	PARIS (Reuters) - French President Emmanuel Ma...	worldnews	August 31, 2017
21263	Guatemala top court sides with U.N. graft unit...	GUATEMALA CITY (Reuters) - Guatemala s top cou...	worldnews	August 29, 2017
21290	Europeans, Africans agree renewed push to tack...	PARIS (Reuters) - Europe s big four continen...	worldnews	August 28, 2017
21353	Thailand's ousted PM Yingluck has fled abroad:...	BANGKOK (Reuters) - Ousted Thai prime minister...	worldnews	August 25, 2017
21408	U.S., North Korea clash at U.N. forum over nuc...	GENEVA (Reuters) - North Korea and the United ...	worldnews	August 22, 2017

206 rows x 4 columns

```
# Remove duplicate rows and keep the first occurrence
df = df.drop_duplicates()
```

```
import string
```

```
# Function to remove punctuation
def remove_punctuation(text):
    if isinstance(text, str):
        return text.translate(str.maketrans('', '', string.punctuation))
    return text
```

```
# Apply to all object (text) columns
for col in df.select_dtypes(include='object').columns:
    df[col] = df[col].apply(remove_punctuation)
```

MERGE:(ONE-HOT ENCODING)

```
import pandas as pd
```

```
# Load CSVs with encoding
```

```
df_true = pd.read_csv('True.csv', encoding='utf-8')
df_fake = pd.read_csv('Fake.csv', encoding='utf-8')

# Add label column
df_true['label'] = 'REAL'
df_fake['label'] = 'FAKE'

# Merge DataFrames
df = pd.concat([df_true, df_fake], ignore_index=True)

# One-hot encode the 'label' column
df_encoded = pd.get_dummies(df, columns=['label'])

# Save to new CSV file with utf-8 encoding
df_encoded.to_csv('News.csv', index=False, encoding='utf-8')

# Optional: show confirmation
print("News.csv saved with one-hot encoding and utf-8 encoding.")
```

News.csv saved with one-hot encoding and utf-8 encoding.

ENCODING:

```
import pandas as pd

# Load both files with UTF-8 encoding
df_true = pd.read_csv('True.csv', encoding='utf-8')
df_fake = pd.read_csv('Fake.csv', encoding='utf-8')

# Add label columns
df_true['label'] = 'REAL'
df_fake['label'] = 'FAKE'

# Merge both DataFrames
df_merged = pd.concat([df_true, df_fake], ignore_index=True)

# Save merged DataFrame with UTF-8 encoding
df_merged.to_csv('News.csv', index=False, encoding='utf-8')

print("News.csv created with UTF-8 encoding.")
```

News.csv created with UTF-8 encoding.

VISUALIZATION:

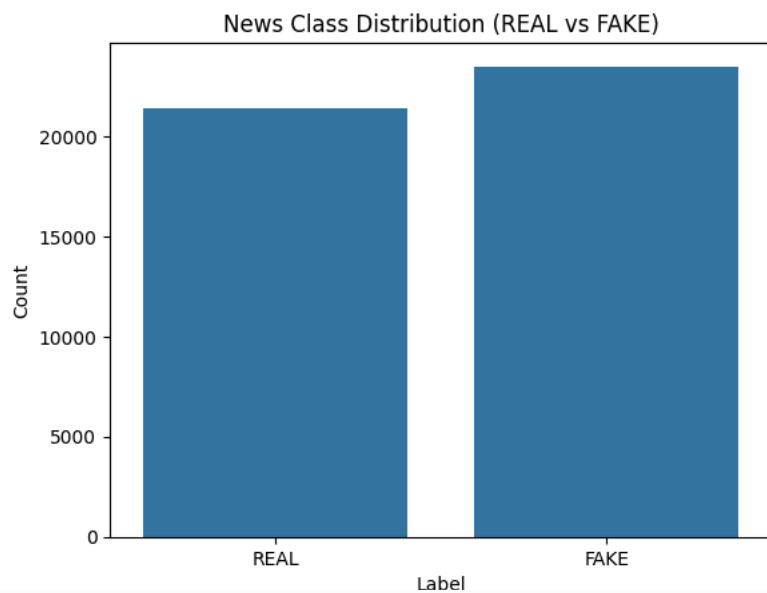
1. BAR PLOT

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

# Load the merged and one-hot encoded dataset
df = pd.read_csv('News.csv', encoding='utf-8')

# --- If you still have 'label' column as categorical (REAL/FAKE) ---
# If not, re-create from one-hot columns
if 'label' not in df.columns and 'label_REAL' in df.columns:
    df['label'] = df['label_REAL'].apply(lambda x: 'REAL' if x == 1 else 'FAKE')

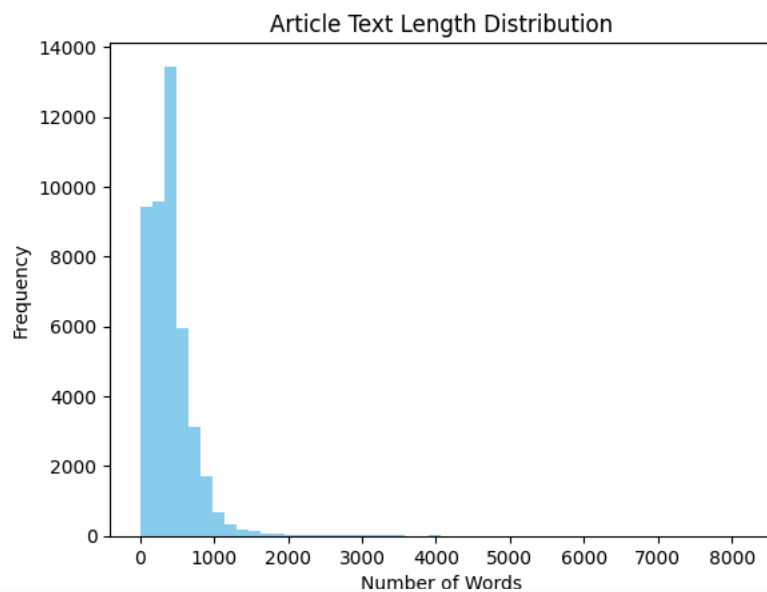
# Plot label distribution
sns.countplot(data=df, x='label')
plt.title('News Class Distribution (REAL vs FAKE)')
plt.xlabel('Label')
plt.ylabel('Count')
plt.show()
```



2.HISTOGRAM

```
# Add a column for text length
df['text_length'] = df['text'].apply(lambda x: len(str(x).split()))

# Plot histogram
plt.hist(df['text_length'], bins=50, color='skyblue')
plt.title('Article Text Length Distribution')
plt.xlabel('Number of Words')
plt.ylabel('Frequency')
plt.show()
```



COMPLETE ONE HOT ENCODING:

```
import pandas as pd

# Step 1: Load the datasets with encoding
df_true = pd.read_csv('True.csv', encoding='utf-8')
df_fake = pd.read_csv('Fake.csv', encoding='utf-8')

# Step 2: Add a 'label' column
df_true['label'] = 'REAL'
df_fake['label'] = 'FAKE'

# Step 3: Merge the datasets
df = pd.concat([df_true, df_fake], ignore_index=True)

# Step 4: One-hot encode the 'label' column
df_encoded = pd.get_dummies(df, columns=['label'])
```

```
# Step 5: Save to new CSV
df_encoded.to_csv('News.csv', index=False, encoding='utf-8')
```

```
# Step 6: Show result
print("One-hot encoded DataFrame saved as 'News.csv'.")
print(df_encoded[['label_FAKE', 'label_REAL']].head())
```

```
↗ One-hot encoded DataFrame saved as 'News.csv'.
  label_FAKE  label_REAL
0         False         True
1         False         True
2         False         True
3         False         True
4         False         True
```

CHECK NULL:

```
import pandas as pd
```

```
# Load the dataset
df = pd.read_csv('News.csv', encoding='utf-8')
```

```
# Check for null values in each column
null_counts = df.isnull().sum()
```

```
# Display columns with missing values
print("Null values per column:")
print(null_counts[null_counts > 0])
df_cleaned = df.dropna()
df_filled = df.fillna('')
```

```
↗ Null values per column:
Series([], dtype: int64)
```

TRAIN TEST:

```
import pandas as pd
from sklearn.model_selection import train_test_split
```

```
# Load dataset
df = pd.read_csv('News.csv', encoding='utf-8')
```

```
# Choose features and target
X = df['text'] # or df[['title', 'text']] if using multiple columns
y = df[['label_FAKE', 'label_REAL']] # One-hot encoded target
```

```
# Split into train/test (e.g., 80% train, 20% test)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

```
# Check sizes
print(f"Training samples: {len(X_train)}")
print(f"Testing samples: {len(X_test)}")
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, stratify=y, random_state=42)
```

```
↗ Training samples: 35918
Testing samples: 8980
```

TEXT LENGTH FEATURE:

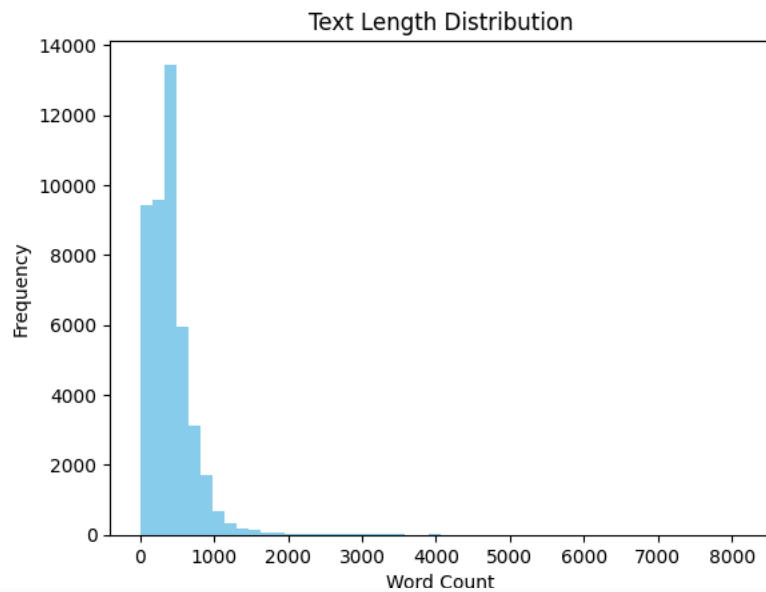
```
df['text_length'] = df['text'].apply(lambda x: len(str(x).split()))
print(df['text_length'].describe())
```

```
# Histogram of article length
df['text_length'].plot.hist(bins=50, title='Text Length Distribution', color='skyblue')
plt.xlabel('Word Count')
plt.show()
```

```

count    44898.000000
mean      405.282284
std       351.265595
min        0.000000
25%       203.000000
50%       362.000000
75%       513.000000
max      8135.000000
Name: text_length, dtype: float64

```



MODEL BUILDING:

```

import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
import matplotlib.pyplot as plt
import seaborn as sns

# Step 1: Load dataset
df = pd.read_csv('News.csv', encoding='utf-8')

# Step 2: Prepare data
df['label'] = df['label_REAL'].apply(lambda x: 1 if x == 1 else 0) # 1 = REAL, 0 = FAKE
X = df['text']
y = df['label']

# Step 3: Train/test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42, stratify=y)

# Step 4: TF-IDF vectorization
tfidf = TfidfVectorizer(stop_words='english', max_df=0.7)
X_train_vec = tfidf.fit_transform(X_train)
X_test_vec = tfidf.transform(X_test)

# Step 5: Train logistic regression model
model = LogisticRegression()
model.fit(X_train_vec, y_train)

# Step 6: Predictions and evaluation
y_pred = model.predict(X_test_vec)

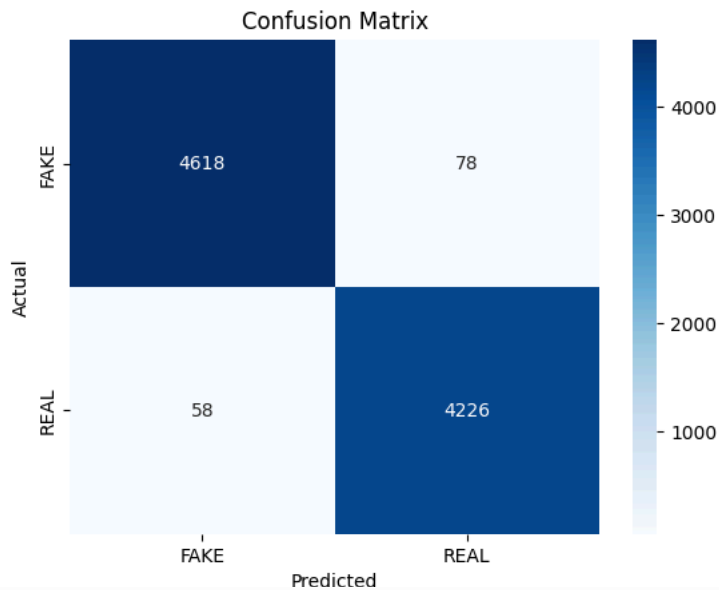
print("Accuracy:", accuracy_score(y_test, y_pred))
print("\nClassification Report:\n", classification_report(y_test, y_pred))

# Confusion Matrix
cm = confusion_matrix(y_test, y_pred)
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=['FAKE', 'REAL'], yticklabels=['FAKE', 'REAL'])
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.title('Confusion Matrix')
plt.show()

```


Accuracy: 0.9848552338530067

Classification Report:				
	precision	recall	f1-score	support
0	0.99	0.98	0.99	4696
1	0.98	0.99	0.98	4284
accuracy			0.98	8980
macro avg	0.98	0.98	0.98	8980
weighted avg	0.98	0.98	0.98	8980



MODEL EVALUATION:

```
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
import seaborn as sns
import matplotlib.pyplot as plt

# Load data
df = pd.read_csv('News.csv', encoding='utf-8')
df['label'] = df['label_REAL'].apply(lambda x: 1 if x == 1 else 0) # 1 = REAL, 0 = FAKE

# Split
X_train, X_test, y_train, y_test = train_test_split(df['text'], df['label'], test_size=0.2, random_state=42, stratify=df['label'])

# TF-IDF
tfidf = TfidfVectorizer(stop_words='english', max_df=0.7)
X_train_vec = tfidf.fit_transform(X_train)
X_test_vec = tfidf.transform(X_test)

# Model
model = LogisticRegression()
model.fit(X_train_vec, y_train)
y_pred = model.predict(X_test_vec)

# Evaluation
print("Accuracy:", accuracy_score(y_test, y_pred))
print("\nClassification Report:\n", classification_report(y_test, y_pred))

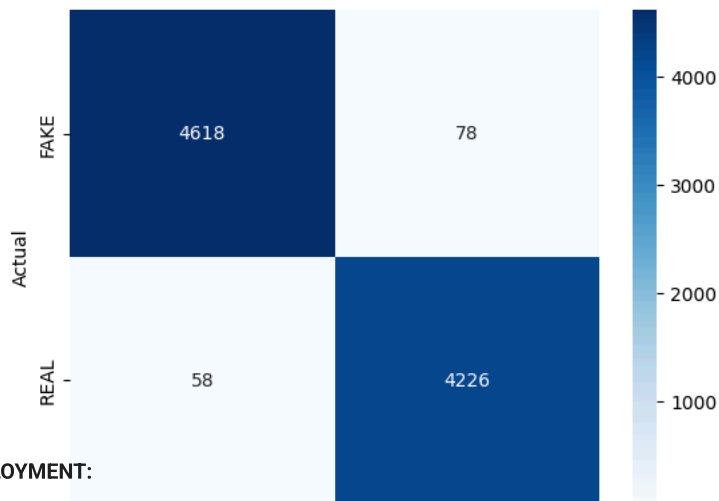
# Confusion matrix
cm = confusion_matrix(y_test, y_pred)
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=['FAKE', 'REAL'], yticklabels=['FAKE', 'REAL'])
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.title('Confusion Matrix')
plt.show()
```

Accuracy: 0.9848552338530067

Classification Report:

	precision	recall	f1-score	support
0	0.99	0.98	0.99	4696
1	0.98	0.99	0.98	4284
accuracy			0.98	8980
macro avg	0.98	0.98	0.98	8980
weighted avg	0.98	0.98	0.98	8980

Confusion Matrix



```
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
import joblib

# Step 1: Load and prepare data
df = pd.read_csv('News.csv', encoding='utf-8')
df['label'] = df['label_REAL'].apply(lambda x: 1 if x == 1 else 0) # 1 = REAL, 0 = FAKE
X_train, X_test, y_train, y_test = train_test_split(df['text'], df['label'], test_size=0.2, random_state=42, stratify=df['label'])

# Step 2: Vectorization
tfidf = TfidfVectorizer(stop_words='english', max_df=0.7)
X_train_vec = tfidf.fit_transform(X_train)
X_test_vec = tfidf.transform(X_test)

# Step 3: Train model
model = LogisticRegression()
model.fit(X_train_vec, y_train)
```