


```

1 import pandas as pd
2 import numpy as np
3 import re
4 import string
5 import nltk
6 import matplotlib.pyplot as plt
7 import seaborn as sns
8
9 from nltk.corpus import stopwords
10 from sklearn.feature_extraction.text import TfidfVectorizer
11 from sklearn.model_selection import train_test_split
12 from sklearn.svm import LinearSVC
13 from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
14
15 nltk.download('stopwords')
16 stop_words = set(stopwords.words('english'))


```

 [nltk\_data] Downloading package stopwords to /root/nltk\_data...  
[nltk\_data] Unzipping corpora/stopwords.zip.

```

1 df = pd.read_csv("news_dataset.csv") # Make sure this file exists
2 print("Dataset shape:", df.shape)
3 df = df[['text', 'label']].dropna()
4


```

 Dataset shape: (3729, 2)

```

1 df.head()

```



	text	label
0	Payal has accused filmmaker Anurag Kashyap of ...	REAL
1	A four-minute-long video of a woman criticisin...	FAKE
2	Republic Poll, a fake Twitter account imitatin...	FAKE
3	Delhi teen finds place on UN green list, turns...	REAL
4	Delhi: A high-level meeting underway at reside...	REAL

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

```

1 def preprocess_text(text):
2     text = str(text).lower() # lowercase
3     text = re.sub(r"http\S+|www\S+|https\S+", '', text) # remove links
4     text = re.sub(r"^\w\s", '', text) # remove punctuation
5     text = re.sub(r"\d+", '', text) # remove numbers

```

```

6 text = ' '.join([word for word in text.split() if word not in stop_words])
7 return text
8

```

```

1 df['cleaned_text'] = df['text'].apply(preprocess_text)
2 print("\nSample cleaned text:\n", df['cleaned_text'].head())

```



Sample cleaned text:

```

0    payal accused filmmaker anurag kashyap behavin...
1    fourminutelong video woman criticising governm...
2    republic poll fake twitter account imitating a...
3    delhi teen finds place un green list turns gla...
4    delhi highlevel meeting underway residence raj...
Name: cleaned_text, dtype: object

```

```

1 tfidf = TfidfVectorizer(max_features=5000)
2 X = tfidf.fit_transform(df['cleaned_text'])
3 y = df['label']

```

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

```

1 model = LinearSVC()
2 model.fit(X_train, y_train)

```



LinearSVC ⓘ ?

LinearSVC()

```

1 y_pred = model.predict(X_test)
2 accuracy = accuracy_score(y_test, y_pred)
3 print("\n model accuracy:", accuracy)
4 print("\n classification report:\n", classification_report(y_test, y_pred))

```



model accuracy: 0.9986577181208054

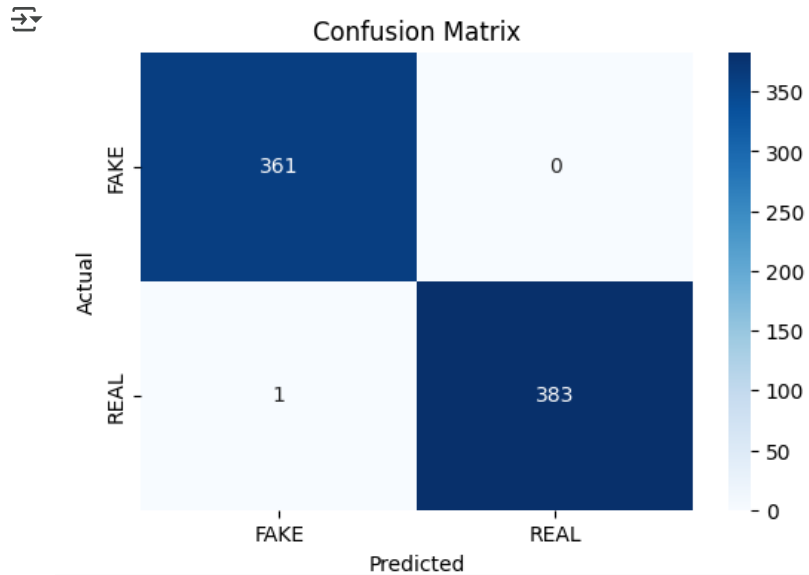
classification report:

	precision	recall	f1-score	support
FAKE	1.00	1.00	1.00	361
REAL	1.00	1.00	1.00	384
accuracy			1.00	745
macro avg	1.00	1.00	1.00	745
weighted avg	1.00	1.00	1.00	745

```

1 cm = confusion_matrix(y_test, y_pred)
2 plt.figure(figsize=(6,4))
3 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=['FAKE', 'REAL'], yticklabels=['FAKE', 'REAL'])
4 plt.title('Confusion Matrix')
5 plt.xlabel('Predicted')
6 plt.ylabel('Actual')
7 plt.show()

```



```

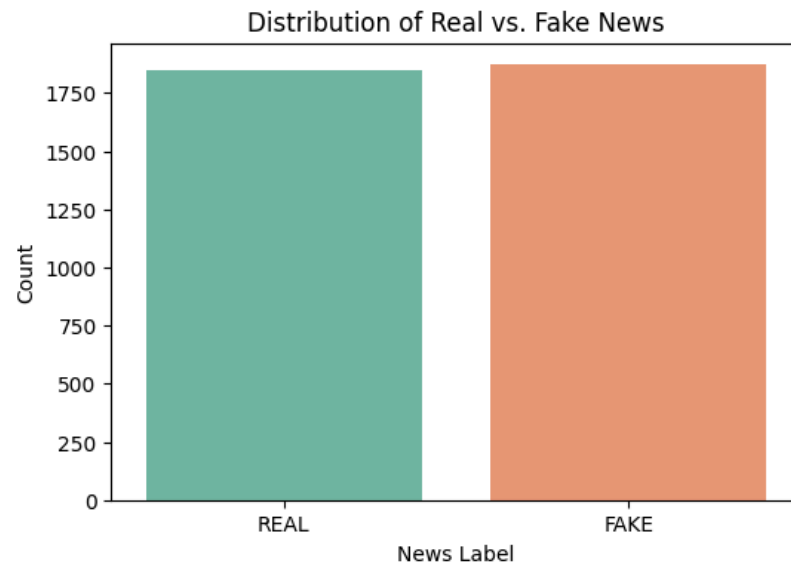
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 # Count plot
5 plt.figure(figsize=(6,4))
6 sns.countplot(x='label', data=df, palette='Set2')
7 plt.title('Distribution of Real vs. Fake News')
8 plt.xlabel('News Label')
9 plt.ylabel('Count')
10 plt.show()
11

```

 /tmp/ipython-input-11-1299514214.py:6: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the same effect.

```
sns.countplot(x='label', data=df, palette='Set2')
```



```
1 # Pie chart
2 label_counts = df['label'].value_counts()
3 plt.figure(figsize=(5,5))
4 plt.pie(label_counts, labels=label_counts.index, autopct='%1.1f%%', startangle=140, colors=['#66bb6a', '#ef5350'])
5 plt.title('Percentage of Real vs. Fake News')
6 plt.axis('equal')
7 plt.show()
8
```



Percentage of Real vs. Fake News

