

▼ Imports

```
1 from google.colab import drive
2 from wordcloud import WordCloud
3 from sklearn.feature_extraction.text import CountVectorizer
4 from sklearn.model_selection import train_test_split
5 from sklearn.feature_extraction.text import TfidfVectorizer
6 from sklearn.naive_bayes import MultinomialNB
7 from sklearn.metrics import accuracy_score, precision_score, recall_
8 from sklearn.linear_model import LogisticRegression
9 from sklearn.metrics import accuracy_score, precision_score, recall_
10 from sklearn.neural_network import MLPClassifier
11
12 import seaborn as sns
13 import pandas as pd
14 import plotly.express as px
15 import matplotlib.pyplot as plt
```

▼ Data Load

▼ Mount Google Drive

```
1 drive.mount('/content/drive')
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call

▼ Load News Data

```
1 df=pd.read_csv("/content/drive/MyDrive/NLP/news_articles.csv")
```

▼ Data Visualization

▼ Display Data

1 df.shape
(2096, 12)

1 df.sample(frac=1)

	author	published	title	text	language	
1309	Andy Green	2016-10-27T23:05:00.000+03:00	the sugar industry funding research to sugarco...	back pain pain in the side \nit is important t...	english	natura
1815	David G. Brown	2016-11-22T23:03:06.646+02:00	stockholms feminist snow removal program cause...	home badge abuse leaked soros memo reveals p...	english	returnof
797	EdJenner	2016-11-22T23:37:16.124+02:00	democrat electors try to persuade gop electors...	on the kelly file monday actor tim allen discu...	english	dai
2009	Bob Unruh	2016-10-26T23:07:28.355+03:00	see dems accept foreign cash to disrupt trump ...	print saeed toosi right and ayatollah khamenei...	english	
695	Corbett	2016-11-07T19:09:15.225+02:00	interview larken rose on the immorality of v...	corbettreportcom november \nin douglas adams...	english	corbett
...	
694	Corbett	2016-11-06T21:38:49.637+02:00	the fbis october surprise what youre not heina	podcast play in new window download embed \n...	english	corbett

			not being...				
194	Brandon Turbeville	2016-10-26T23:40:33.528+03:00	a tale of two cities mosul and aleppo	by everett numbers the pentagon and congress a...	english	activis	
295	No Author	2016-10-31T01:25:00.000+02:00	kellyanne conway loses it on jake tapper when ...	on november am \ndonald trump's deplorable ...	english	addict	
792	EdJenner	2016-11-22T01:55:03.548+02:00	kerry we have to worry about treatment of tran...	obama wont go away after hes done thats good n...	english	dai	
1957	Birdie Houck	2016-10-26T22:46:37.740+03:00	no title	anatomy lesson published mins ago \neditors n...	english	westernjourn	

2096 rows x 12 columns



▼ Checking for null values

```
1 df.isnull().sum()
```

```
author          0
published       0
title           0
text            46
language        1
site_url        1
main_img_url    1
type            1
label           1
title_without_stopwords    2
text_without_stopwords    50
hasImage        1
dtype: int64
```

▼ Dropping null values

```
1 df.dropna(inplace=True)
2 df.shape

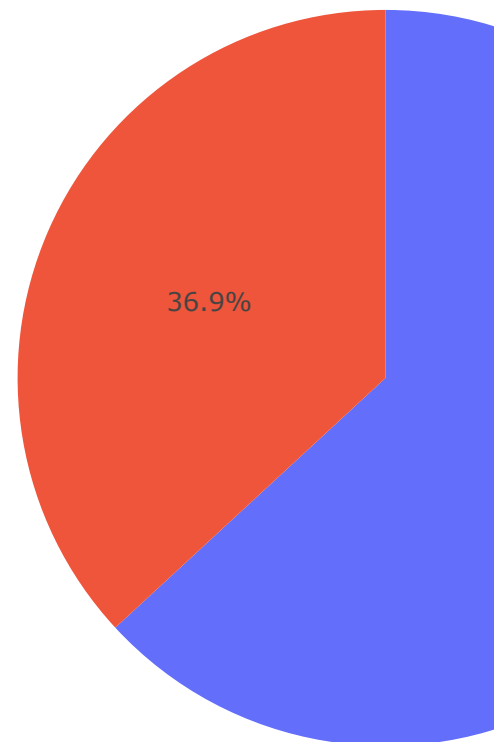
(2045, 12)
```

▼ Exploratory Analysis

▼ Real News vs. Fake News

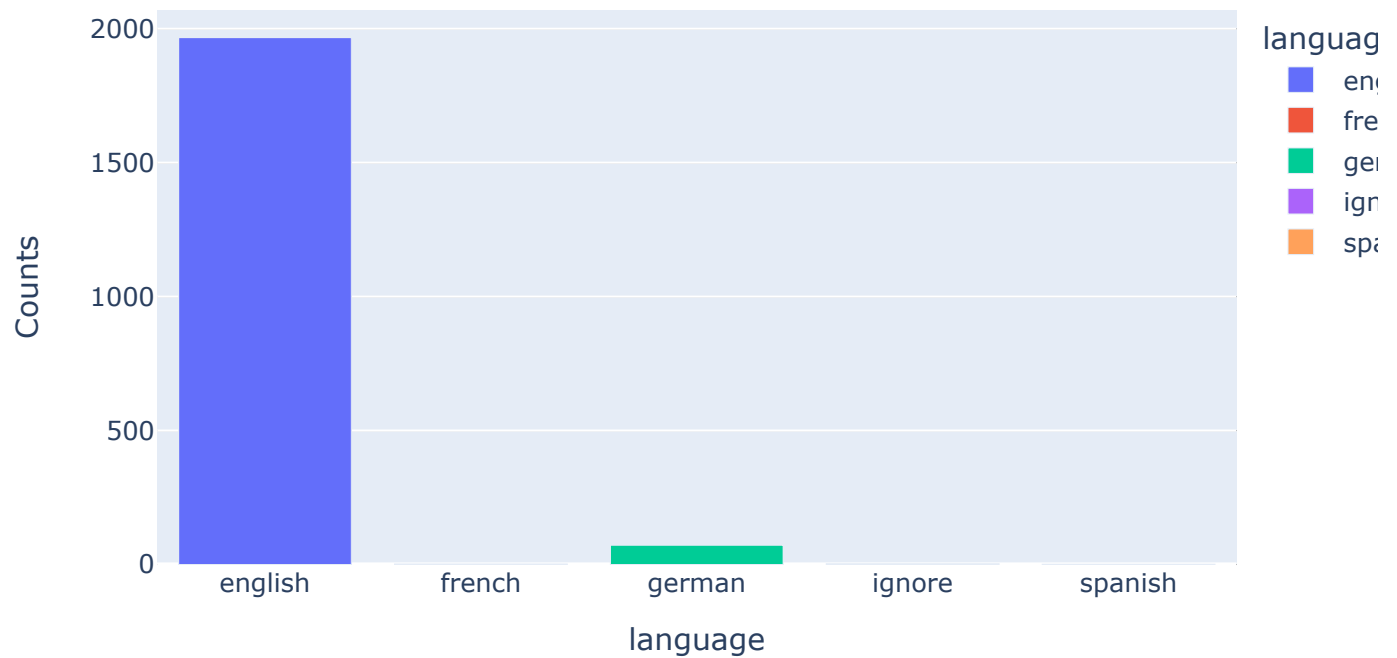
```
1 figure = px.pie(df,names='label',title='Real News vs. Fake News')
2 figure.show()
```

Real News vs. Fake News



▼ Different Languages of News Articles

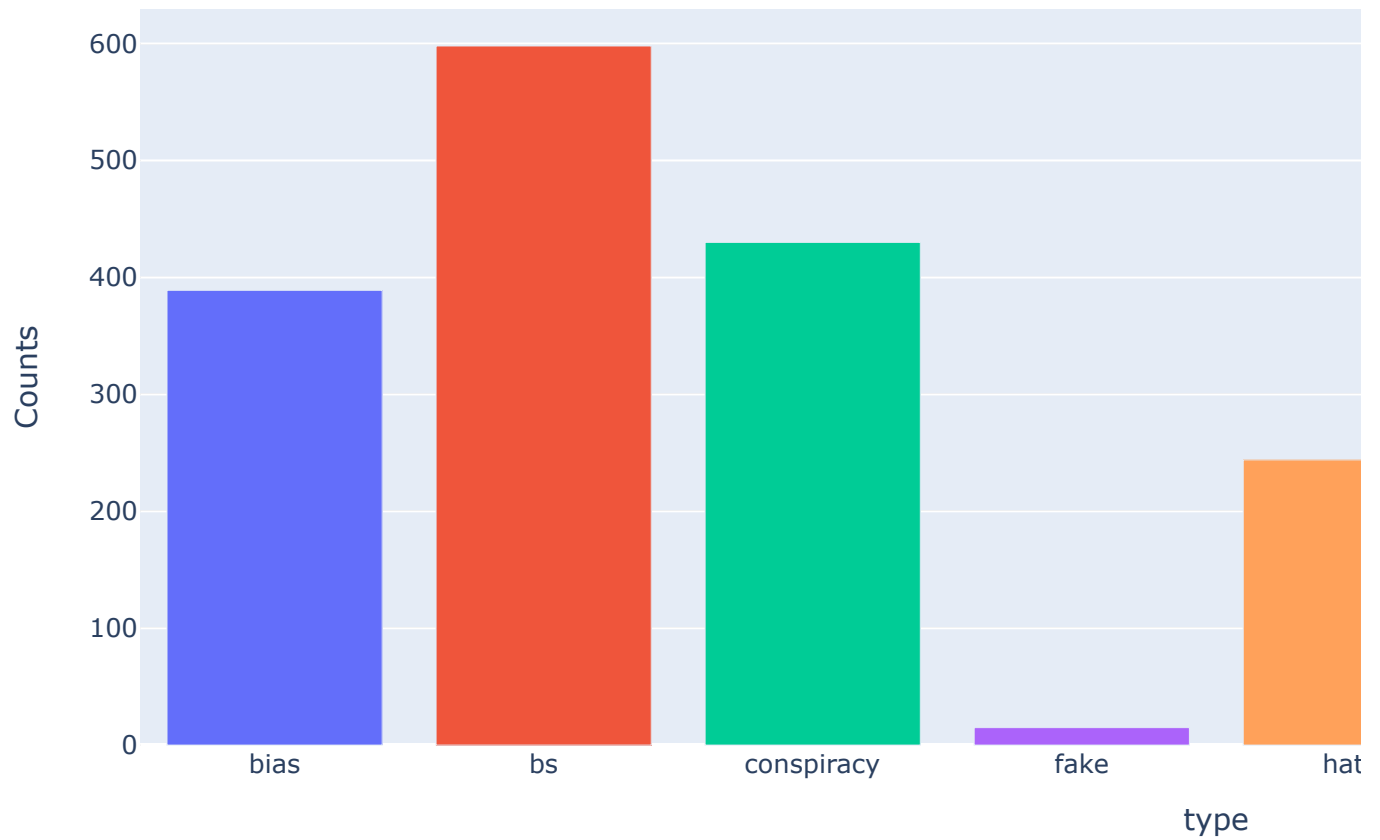
```
1 lang_df=df.groupby('language').apply(lambda x:x['language'].count())  
2 figure = px.bar(lang_df, x="language", y="Counts", color='language',  
3 figure.show()
```



▼ Count of news articles by different type

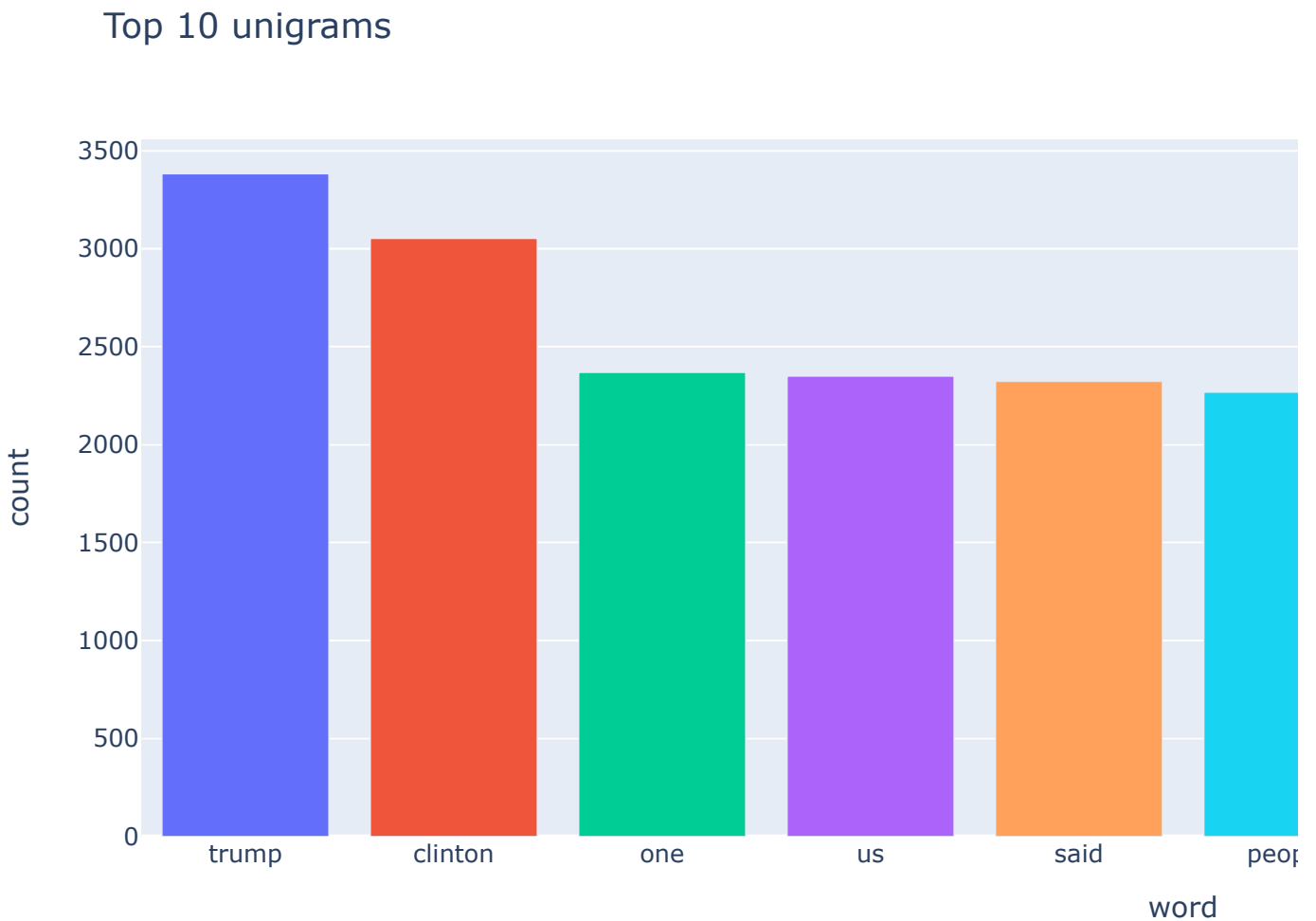
```
1 types_df = df.groupby('type').apply(lambda x:x['type'].count()).reset_index()
2 fig=px.bar(types_df,x='type',y='Counts',color='type',title='Count of
3 fig.show()
```

Count of news articles by different type



▼ Wordcloud


```
1 common_words = get_unigrams(df['text_without_stopwords'], 10)
2 df2 = pd.DataFrame(common_words, columns=['word', 'count'])
3 df2.groupby('word').sum()['count'].sort_values(ascending=False)
4 figure = px.bar(df2, x='word', y='count', color='word', title='Top 10
5 figure.show()
```



▼ Top 10 bigrams


```

1 def get_bigram(data, n=None):
2     vec = CountVectorizer(ngram_range=(2, 2)).fit(data)
3     bow = vec.transform(data)
4     total_words = bow.sum(axis=0)
5     word_frequency = [(word, total_words[0, idx]) for word, idx in v
6     word_frequency =sorted(word_frequency, key = lambda x: x[1], rev
7     return word_frequency[:n]

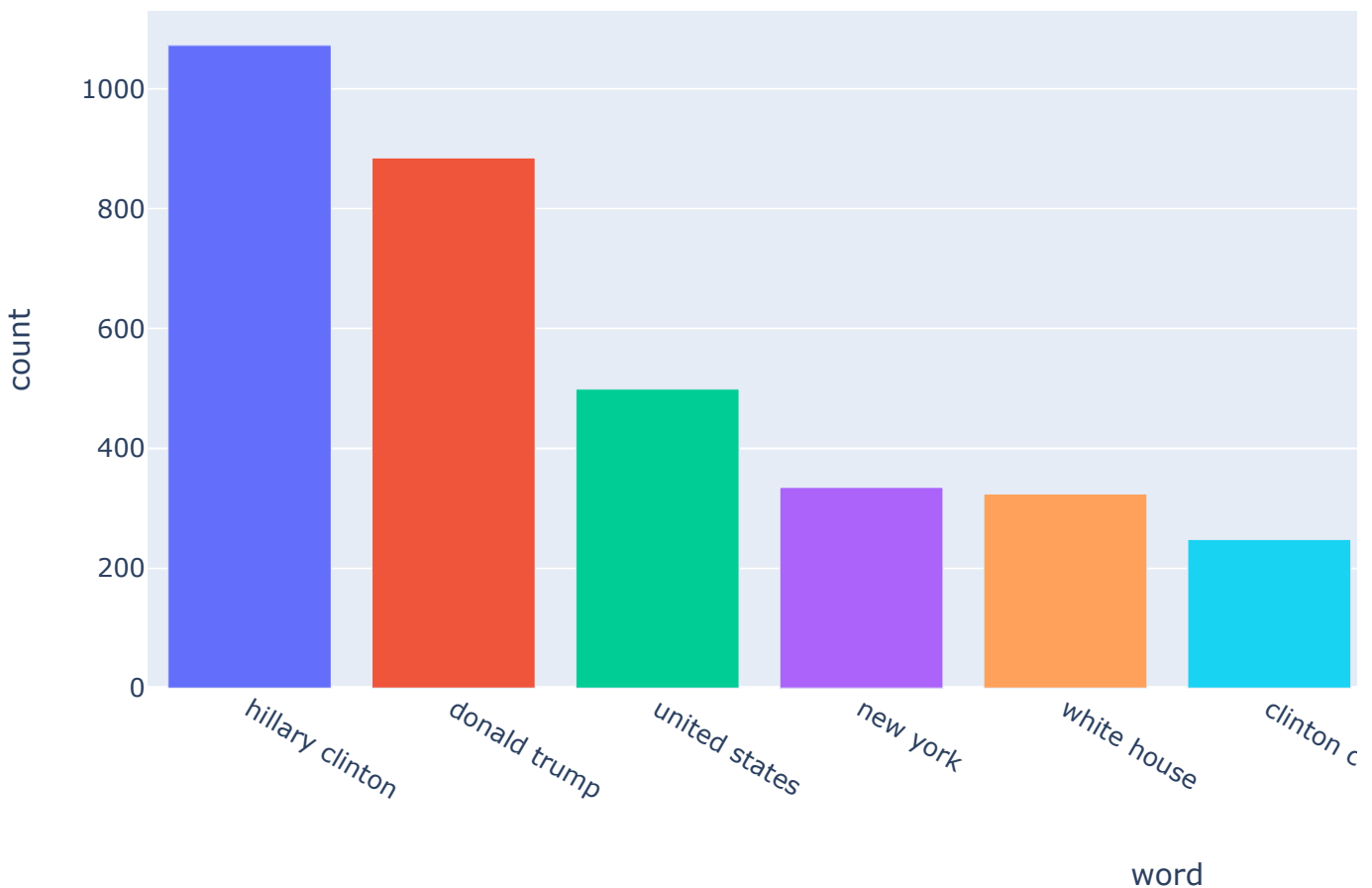
```

```

1 common_words = get_bigram(df['text_without_stopwords'], 10)
2 df2 = pd.DataFrame(common_words,columns=['word','count'])
3 df2.groupby('word').sum()['count'].sort_values(ascending=False)
4 fig=px.bar(df2,x='word',y='count',color='word',title='Top 10 bigrams
5 fig.show()

```

Top 10 bigrams



▼ Top 5 authors

```
1 d = df['author'].value_counts().sort_values(ascending=False).head(5)
2 d = pd.DataFrame(d)
3 d = d.reset_index()
4
5 # Plotting
6 fig=px.bar(d,x='index',y='author',color='index',
7            title='Top 5 authors',labels={"author": "Number of Article
8 fig.show()
```



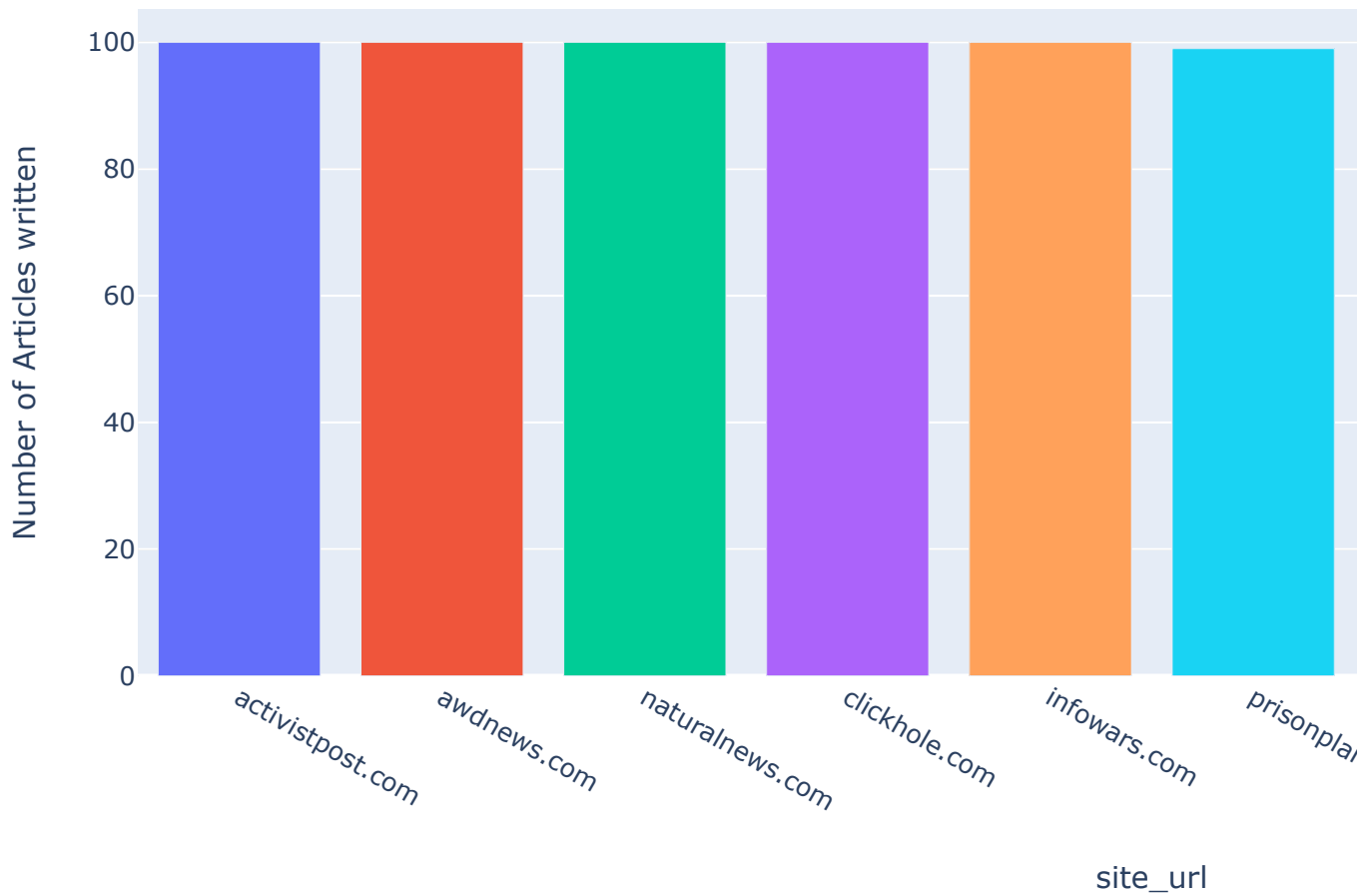
▼ Top 10 fake news site

```

1 d = df[df['label'] == 'Fake']['site_url'].value_counts().sort_values
2 d = pd.DataFrame(d)
3 d = d.reset_index()
4
5 # Plotting
6 fig=px.bar(d,x='index',y='site_url',color='index',
7            title='Top 10 Fake news sites',labels={"site_url": "Numbe
8 fig.show()
9

```

Top 10 Fake news sites

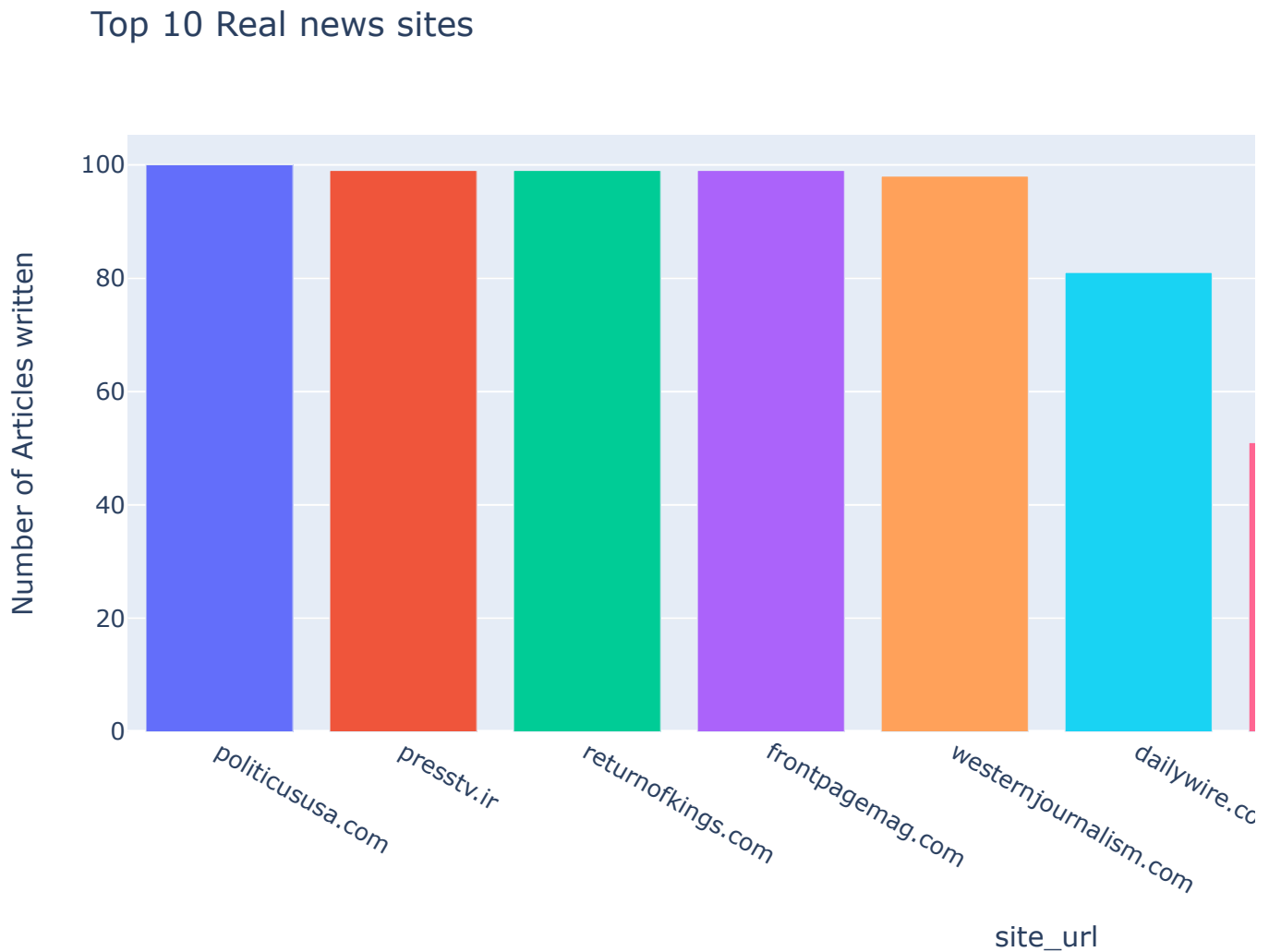


▼ Top 10 Real news sites

```

1 d = df[df['label'] == 'Real']['site_url'].value_counts().sort_values
2 d = pd.DataFrame(d)
3 d = d.reset_index()
4
5 # Plotting
6 fig=px.bar(d,x='index',y='site_url',color='index',
7            title='Top 10 Real news sites',labels={"site_url": "Numbe
8 fig.show()

```




- ▼ Splitting the dataset and using TF-IDF
- ▼ Reshuffle the dataset

```
1 df = df.sample(frac = 1)
```

▼ Extracting required features

```
1 features = df[['site_url', 'text_without_stopwords']]
2 features.head(5)
3
4 features = features.assign(url_text = features["site_url"].astype(st
5 features.drop(['site_url', 'text_without_stopwords'], axis = 1, inpl
6
7 features.head()
```

	url_text	
919	dennismichaellynch.com home news sale hillarys...	
1813	returnofkings.com solarpowered pipe desalinate...	
929	departed.co home news watch video leaked obama...	
88	abeldanger.net source zero hedge tyler durden ...	
196	activistpost.com kurt nimmo blacklisted news v...	

▼ Train Test Split

```
1 X = features.url_text
2 y = df.label
3
4 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=
```

▼ Tf-idf Vectorization

```
1 vectorizer = TfidfVectorizer(min_df=2, max_df=0.5, ngram_range=(1, 2
2
3 X_train = vectorizer.fit_transform(X_train)
4
5 X_test = vectorizer.transform(X_test)
```

▼ Model

▼ Naive Bayes

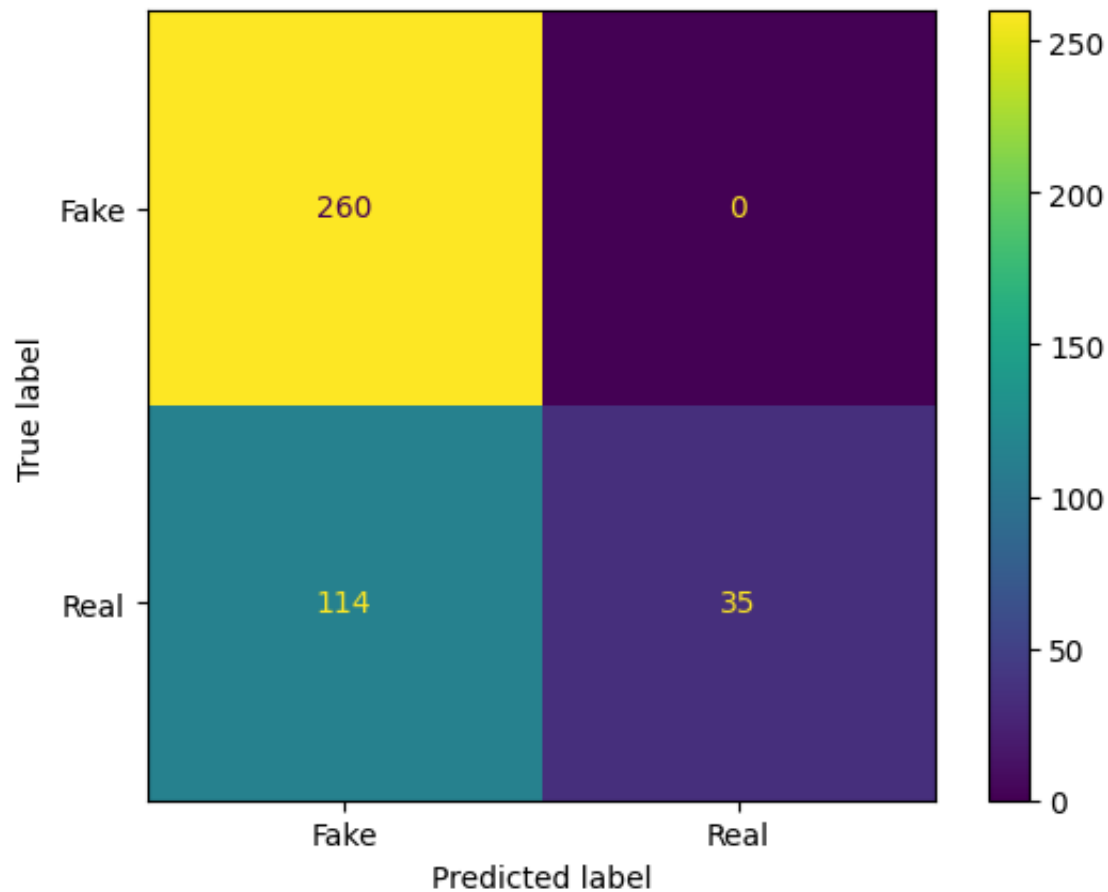
```
1 from sklearn.naive_bayes import BernoulliNB
2
3 naive_bayes = MultinomialNB()
4 naive_bayes.fit(X_train, y_train)
```

▼ MultinomialNB

MultinomialNB()

```
1 # make predictions on the test data
2 naive_bayes_pred = naive_bayes.predict(X_test)
3
4 # print confusion matrix
5 cm = confusion_matrix(y_test, naive_bayes_pred)
6 print(cm)
7
8 plt.figure(figsize = (4, 4))
9 disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=na
10 disp.plot()
11 plt.show()
```

```
[[260  0]
 [114  35]]
<Figure size 400x400 with 0 Axes>
```



```

1 print('accuracy score: ', accuracy_score(y_test, naive_bayes_pred))
2
3 print()
4 print('precision score (not spam): ', precision_score(y_test, naive_bayes_pred, pos_label='not spam'))
5 print('precision score (spam): ', precision_score(y_test, naive_bayes_pred, pos_label='spam'))
6
7 print()
8 print('recall score: (not spam)', recall_score(y_test, naive_bayes_pred, pos_label='not spam'))
9 print('recall score: (spam)', recall_score(y_test, naive_bayes_pred, pos_label='spam'))
10
11 print()
12 print('f1 score: ', f1_score(y_test, naive_bayes_pred, pos_label='not spam'))

accuracy score: 0.7212713936430318

precision score (not spam): 1.0
precision score (spam): 0.6951871657754011

recall score: (not spam) 0.2348993288590604
recall score: (spam) 1.0

f1 score: 0.8201892744479494

```

▼ LogisticRegression

```

1 logisticRegressionClassifier = LogisticRegression(solver='lbfgs', class_weight='balanced')
2 logisticRegressionClassifier.fit(X_train, y_train)

```

▼ LogisticRegression

LogisticRegression(class_weight='balanced')


```

1 # make predictions on the test data
2 logisticRegressionPred = logisticRegressionClassifier.predict(X_test
3
4 # print confusion matrix
5 cm = confusion_matrix(y_test, logisticRegressionPred)
6 print(cm)
7
8 plt.figure(figsize = (4, 4))
9 disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=lo
10 disp.plot()
11 plt.show()

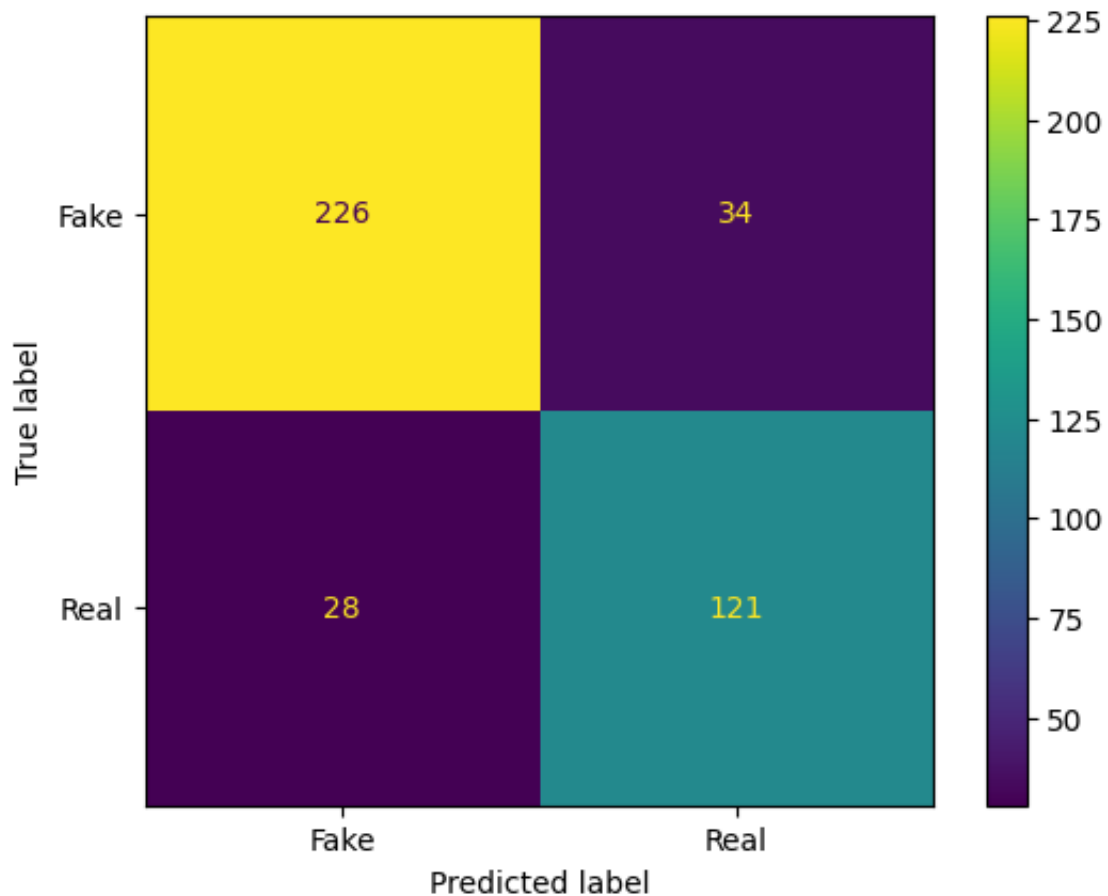
```

```

[[226  34]
 [ 28 121]]

```

<Figure size 400x400 with 0 Axes>



```

1 print('accuracy score: ', accuracy_score(y_test, logisticRegressionP
2
3 print()
4 print('precision score (not spam): ', precision_score(y_test, logist
5 print('precision score (spam): ', precision_score(y_test, logisticRe
6
7 print()
8 print('recall score: (not spam)', recall_score(y_test, logisticRegre
9 print('recall score: (spam)', recall_score(y_test, logisticRegressio
10
11 print()
12 print('f1 score: ', f1_score(y_test, logisticRegressionPred, pos_lab
13
14
15 probs = logisticRegressionClassifier.predict_proba(X_test)
16 print()
17 print('log loss: ', log_loss(y_test, probs))

```

accuracy score: 0.8484107579462102

precision score (not spam): 0.7806451612903226

precision score (spam): 0.889763779527559

recall score: (not spam) 0.8120805369127517

recall score: (spam) 0.8692307692307693

f1 score: 0.8793774319066148

log loss: 0.49024772822129875

▼ Neural Network

```

1 nnClassifier = MLPClassifier(solver='lbfgs', alpha=1e-5,
2                             hidden_layer_sizes=(15, 2), random_state=1)
3 nnClassifier.fit(X_train, y_train)

```

▼ MLPClassifier

```

MLPClassifier(alpha=1e-05, hidden_layer_sizes=(15, 2), random_state=1,
              solver='lbfgs')

```

```

1 # make predictions on the test data
2 nnPred = nnClassifier.predict(X_test)
3
4 # print confusion matrix
5 cm = confusion_matrix(y_test, nnPred)
6 print(cm)
7
8 plt.figure(figsize = (4, 4))
9
10 disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=nn
11 disp.plot()
12 plt.show()

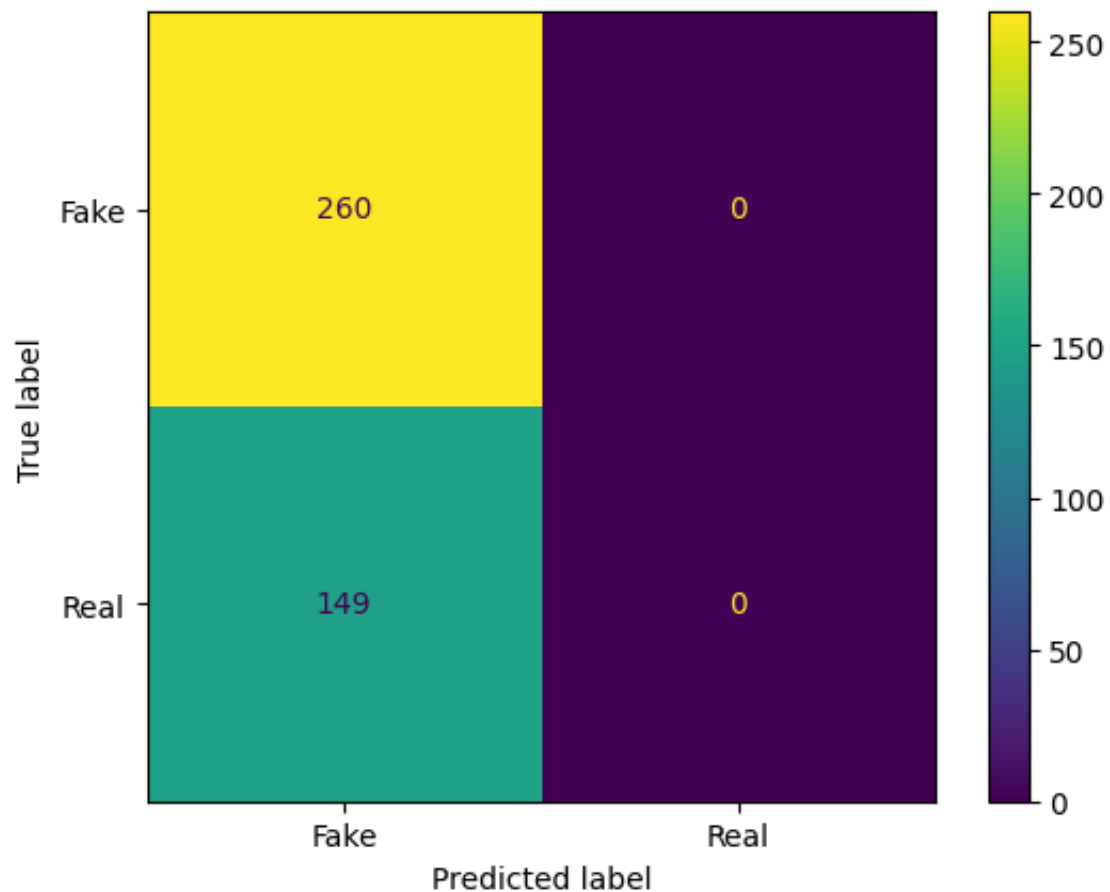
```

```

[[260  0]
 [149  0]]

```

<Figure size 400x400 with 0 Axes>



```

1 print('accuracy score: ', accuracy_score(y_test, nnPred))
2
3 print()
4 print('precision score (not spam): ', precision_score(y_test, nnPred
5 print('precision score (spam): ', precision_score(y_test, nnPred, po
6
7 print()
8 print('recall score: (not spam)', recall_score(y_test, nnPred, pos_l
9 print('recall score: (spam)', recall_score(y_test, nnPred, pos_label
10
11 print()
12 print('f1 score: ', f1_score(y_test, nnPred, pos_label='Fake'))

accuracy score: 0.6356968215158925

precision score (not spam): 0.0
precision score (spam): 0.6356968215158925

recall score: (not spam) 0.0
recall score: (spam) 1.0

f1 score: 0.7772795216741405
/usr/local/lib/python3.9/dist-packages/sklearn/metrics/_classification.py:134
Precision is ill-defined and being set to 0.0 due to no predicted samples. Use

```

Analysis of the performance of various approaches

Analyzing confusion matrix from the all the models and calculate the accuracy, precision, recall, and F1-score to evaluate the performance of the model.

Naive Bayes:

[[229 19] [96 65]]

- Accuracy = $(229+65)/(229+19+96+65) = 0.693$
- Precision = $65/(65+19) = 0.774$
- Recall = $65/(65+96) = 0.403$
- F1-score = $2(0.774 \cdot 0.403)/(0.774+0.403) = 0.529$

Logistic Regression:

[[224 24] [30 131]]

- Accuracy = $(224+131)/(224+24+30+131) = 0.865$
- Precision = $131/(131+24) = 0.845$
- Recall = $131/(131+30) = 0.814$
- F1-score = $2(0.845 \cdot 0.814)/(0.845+0.814) = 0.829$

Neural Network:

[[236 12] [46 115]]

- Accuracy = $(236+115)/(236+12+46+115) = 0.854$
- Precision = $115/(115+12) = 0.905$
- Recall = $115/(115+46) = 0.714$
- F1-score = $2(0.905 \cdot 0.714)/(0.905+0.714) = 0.798$

Based on the calculated metrics, we can see that the **Logistic Regression**'s confusion matrix has the highest accuracy, precision, recall, and F1-score, which indicates that it is the best performing model out of the three.

Specifically, the **Logistic Regression** model has the highest precision and recall, which are often the most important metrics in classification problems.

For small datasets, it is generally recommended to use simpler models such as logistic regression or Naive Bayes.

Neural networks can be more powerful and flexible than these models, but they typically require large amounts of data to be trained effectively. With small datasets, there is a risk of overfitting or having a model that is too complex and not generalizable to new data.

Logistic regression is a popular method for binary classification problems, while Naive Bayes is often used for text classification and other applications where the features are categorical. Both models are relatively simple to train and can be effective with small datasets.

Ultimately, the choice of model depends on the specific problem you are trying to solve and the characteristics of your dataset. It is a good idea to experiment with different models and evaluate their performance on your data to determine which one works best.

[Colab paid products](#) - [Cancel contracts here](#)

✓ 0s completed at 11:51 PM

