

```

1 from google.colab import files
2 import pandas as pd
3
4 # Upload the file
5 uploaded = files.upload()

```



Choose Files 4 files

- **Sports.csv**(text/csv) - 4974 bytes, last modified: 8/7/2024 - 100% done
 - **Politics.csv**(text/csv) - 5529 bytes, last modified: 8/7/2024 - 100% done
 - **Education.csv**(text/csv) - 5083 bytes, last modified: 8/7/2024 - 100% done
 - **Finance.csv**(text/csv) - 5467 bytes, last modified: 8/7/2024 - 100% done
- Saving Sports.csv to Sports.csv
 Saving Politics.csv to Politics.csv
 Saving Education.csv to Education.csv
 Saving Finance.csv to Finance.csv

```

1 education=pd.read_csv("Education.csv")
2 finance=pd.read_csv("Finance.csv")
3 politics=pd.read_csv("Politics.csv")
4 sports=pd.read_csv("Sports.csv")

1 df = pd.concat([education, finance, politics, sports])
2 df

```



	Text	Label
0	The impact of educational reforms remains unce...	positive
1	Critics argue that recent improvements in the ...	negative
2	Innovative teaching methods have led to unexpe...	positive
3	Despite budget constraints, the school has man...	positive
4	The true effectiveness of online learning plat...	negative
...
51	Sports fandom can foster a sense of community ...	positive
52	Sports events offer a platform for showcasing ...	positive
53	The pressure to win in sports can overshadow t...	negative
54	Sports programs in schools play a crucial role...	positive
55	The commercialization of sports has led to exp...	negative

209 rows x 2 columns

Next steps:

[Generate code with df](#)

[View recommended plots](#)

[New interactive sheet](#)

```
1 df['Label'].unique()
```



```
array(['positive', 'negative'], dtype=object)
```

```

1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer
4 from sklearn.naive_bayes import MultinomialNB, BernoulliNB, GaussianNB
5 from sklearn.metrics import accuracy_score, classification_report
6 import nltk
7 from nltk.corpus import stopwords
8 from nltk.stem import WordNetLemmatizer
9 import string
10 import re
11 import joblib
12 from joblib import load
13 from joblib import dump

1 # Download NLTK data files (only the first time)
2 nltk.download('stopwords')
3 nltk.download('wordnet')

```




```

[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Unzipping corpora/stopwords.zip.
[nltk_data] Downloading package wordnet to /root/nltk_data...
True

```

```
1 # Preprocess the text data
2 def preprocess_text(text):
3     # Convert to lowercase
4     text = text.lower()
5     # Remove punctuation
6     text = text.translate(str.maketrans('', '', string.punctuation))
7     # Remove numbers
8     text = re.sub(r'\d+', '', text)
9     # Remove stopwords
10    stop_words = set(stopwords.words('english'))
11    words = text.split()
12    words = [word for word in words if word not in stop_words]
13    # Lemmatize words
14    lemmatizer = WordNetLemmatizer()
15    words = [lemmatizer.lemmatize(word) for word in words]
16    return ' '.join(words)
17
18 df['cleaned_text'] = df['Text'].apply(preprocess_text)
```

```
1 # Cleaned Text
2 df
```



	Text	Label	cleaned_text
0	The impact of educational reforms remains unce...	positive	impact educational reform remains uncertain de...
1	Critics argue that recent improvements in the ...	negative	critic argue recent improvement school system ...
2	Innovative teaching methods have led to unexpe...	positive	innovative teaching method led unexpected chal...
3	Despite budget constraints, the school has man...	positive	despite budget constraint school managed maint...
4	The true effectiveness of online learning plat...	negative	true effectiveness online learning platform st...
...
51	Sports fandom can foster a sense of community ...	positive	sport fandom foster sense community belonging ...
52	Sports events offer a platform for showcasing ...	positive	sport event offer platform showcasing cultural...
53	The pressure to win in sports can overshadow t...	negative	pressure win sport overshadow enjoyment playin...
54	Sports programs in schools play a crucial role...	positive	sport program school play crucial role charact...
55	The commercialization of sports has led to exp...	negative	commercialization sport led exploitation commo...

209 rows x 3 columns

Next steps:

Generate code with df


 View recommended plots

New interactive sheet

```
1 # Function to apply Gaussian Naive Bayes
2 def apply_gaussian_nb(X_train, X_test, y_train, y_test):
3     gnb = GaussianNB()
4     gnb.fit(X_train.toarray(), y_train)
5     y_pred = gnb.predict(X_test.toarray())
6     accuracy = accuracy_score(y_test, y_pred)
7     report = classification_report(y_test, y_pred)
8     return accuracy, report

1 # TF-IDF Vectorizer
2 tfidf = TfidfVectorizer(max_features=5000)
3 X_tfidf = tfidf.fit_transform(df['cleaned_text'])
4 y = df['Label']
5 X_train_tfidf, X_test_tfidf, y_train, y_test = train_test_split(X_tfidf, y, test_size=0.2, random_state=42)

1 # Multinomial Naive Bayes with TF-IDF
2 nb_tfidf = MultinomialNB()
3 nb_tfidf.fit(X_train_tfidf, y_train)
4 y_pred_tfidf = nb_tfidf.predict(X_test_tfidf)
5 accuracy_tfidf = accuracy_score(y_test, y_pred_tfidf)
6 report_tfidf = classification_report(y_test, y_pred_tfidf)
7 print(f'TF-IDF Vectorizer with MultinomialNB Accuracy: {accuracy_tfidf}')
8 print(f'TF-IDF Vectorizer with MultinomialNB Classification Report:\n{report_tfidf}')
```



TF-IDF Vectorizer with MultinomialNB Accuracy: 0.7619047619047619				
TF-IDF Vectorizer with MultinomialNB Classification Report:				
	precision	recall	f1-score	support
negative	0.87	0.62	0.72	21
positive	0.70	0.90	0.79	21
accuracy			0.76	42
macro avg	0.79	0.76	0.76	42
weighted avg	0.79	0.76	0.76	42

```

1 # Bernoulli Naive Bayes with TF-IDF
2 nb_bernoulli_tfidf = BernoulliNB()
3 nb_bernoulli_tfidf.fit(X_train_tfidf, y_train)
4 y_pred_bernoulli_tfidf = nb_bernoulli_tfidf.predict(X_test_tfidf)
5 accuracy_bernoulli_tfidf = accuracy_score(y_test, y_pred_bernoulli_tfidf)
6 report_bernoulli_tfidf = classification_report(y_test, y_pred_bernoulli_tfidf)
7 print(f'TF-IDF Vectorizer with BernoulliNB Accuracy: {accuracy_bernoulli_tfidf}')
8 print(f'TF-IDF Vectorizer with BernoulliNB Classification Report:\n{report_bernoulli_tfidf}')

```

TF-IDF Vectorizer with BernoulliNB Accuracy: 0.7380952380952381

TF-IDF Vectorizer with BernoulliNB Classification Report:

	precision	recall	f1-score	support
negative	0.81	0.62	0.70	21
positive	0.69	0.86	0.77	21
accuracy			0.74	42
macro avg	0.75	0.74	0.73	42
weighted avg	0.75	0.74	0.73	42

```

1 # Gaussian Naive Bayes with TF-IDF
2 accuracy_gaussian_tfidf, report_gaussian_tfidf = apply_gaussian_nb(X_train_tfidf, X_test_tfidf, y_train, y_test)
3 print(f'TF-IDF Vectorizer with GaussianNB Accuracy: {accuracy_gaussian_tfidf}')
4 print(f'TF-IDF Vectorizer with GaussianNB Classification Report:\n{report_gaussian_tfidf}')

```

TF-IDF Vectorizer with GaussianNB Accuracy: 0.7380952380952381

TF-IDF Vectorizer with GaussianNB Classification Report:

	precision	recall	f1-score	support
negative	0.73	0.76	0.74	21
positive	0.75	0.71	0.73	21
accuracy			0.74	42
macro avg	0.74	0.74	0.74	42
weighted avg	0.74	0.74	0.74	42

```

1 # Count Vectorizer
2 count_vectorizer = CountVectorizer(max_features=5000)
3 X_count = count_vectorizer.fit_transform(df['cleaned_text'])
4 X_train_count, X_test_count, y_train, y_test = train_test_split(X_count, y, test_size=0.2, random_state=42)

```

```

1 # Multinomial Naive Bayes with Count Vectorizer
2 nb_count = MultinomialNB()
3 nb_count.fit(X_train_count, y_train)
4 y_pred_count = nb_count.predict(X_test_count)
5 accuracy_count = accuracy_score(y_test, y_pred_count)
6 report_count = classification_report(y_test, y_pred_count)
7 print(f'Count Vectorizer with MultinomialNB Accuracy: {accuracy_count}')
8 print(f'Count Vectorizer with MultinomialNB Classification Report:\n{report_count}')

```

Count Vectorizer with MultinomialNB Accuracy: 0.6904761904761905

Count Vectorizer with MultinomialNB Classification Report:

	precision	recall	f1-score	support
negative	0.70	0.67	0.68	21
positive	0.68	0.71	0.70	21
accuracy			0.69	42
macro avg	0.69	0.69	0.69	42
weighted avg	0.69	0.69	0.69	42

```

1 # Bernoulli Naive Bayes with Count Vectorizer
2 nb_bernoulli_count = BernoulliNB()
3 nb_bernoulli_count.fit(X_train_count, y_train)
4 y_pred_bernoulli_count = nb_bernoulli_count.predict(X_test_count)
5 accuracy_bernoulli_count = accuracy_score(y_test, y_pred_bernoulli_count)
6 report_bernoulli_count = classification_report(y_test, y_pred_bernoulli_count)
7 print(f'Count Vectorizer with BernoulliNB Accuracy: {accuracy_bernoulli_count}')
8 print(f'Count Vectorizer with BernoulliNB Classification Report:\n{report_bernoulli_count}')

```

Count Vectorizer with BernoulliNB Accuracy: 0.7380952380952381

Count Vectorizer with BernoulliNB Classification Report:

	precision	recall	f1-score	support
negative	0.81	0.62	0.70	21
positive	0.69	0.86	0.77	21
accuracy			0.74	42

macro avg	0.75	0.74	0.73	42
weighted avg	0.75	0.74	0.73	42

```

1 # Gaussian Naive Bayes with Count Vectorizer
2 accuracy_gaussian_count, report_gaussian_count = apply_gaussian_nb(X_train_count, X_test_count, y_train, y_test)
3 print(f'Count Vectorizer with GaussianNB Accuracy: {accuracy_gaussian_count}')
4 print(f'Count Vectorizer with GaussianNB Classification Report:\n{report_gaussian_count}')
```

Count Vectorizer with GaussianNB Accuracy: 0.7142857142857143

Count Vectorizer with GaussianNB Classification Report:

	precision	recall	f1-score	support
negative	0.70	0.76	0.73	21
positive	0.74	0.67	0.70	21
accuracy			0.71	42
macro avg	0.72	0.71	0.71	42
weighted avg	0.72	0.71	0.71	42

```

1 # Binary Vectorizer
2 binary_vectorizer = CountVectorizer(binary=True, max_features=5000)
3 X_binary = binary_vectorizer.fit_transform(df['cleaned_text'])
4 X_train_binary, X_test_binary, y_train, y_test = train_test_split(X_binary, y, test_size=0.2, random_state=42)
```

```

1 # Multinomial Naive Bayes with Binary Vectorizer
2 nb_binary = MultinomialNB()
3 nb_binary.fit(X_train_binary, y_train)
4 y_pred_binary = nb_binary.predict(X_test_binary)
5 accuracy_binary = accuracy_score(y_test, y_pred_binary)
6 report_binary = classification_report(y_test, y_pred_binary)
7 print(f'Binary Vectorizer with MultinomialNB Accuracy: {accuracy_binary}')
8 print(f'Binary Vectorizer with MultinomialNB Classification Report:\n{report_binary}')
```

Binary Vectorizer with MultinomialNB Accuracy: 0.6904761904761905

Binary Vectorizer with MultinomialNB Classification Report:

	precision	recall	f1-score	support
negative	0.70	0.67	0.68	21
positive	0.68	0.71	0.70	21
accuracy			0.69	42
macro avg	0.69	0.69	0.69	42
weighted avg	0.69	0.69	0.69	42

```

1 # Bernoulli Naive Bayes with Binary Vectorizer
2 nb_bernoulli_binary = BernoulliNB()
3 nb_bernoulli_binary.fit(X_train_binary, y_train)
4 y_pred_bernoulli_binary = nb_bernoulli_binary.predict(X_test_binary)
5 accuracy_bernoulli_binary = accuracy_score(y_test, y_pred_bernoulli_binary)
6 report_bernoulli_binary = classification_report(y_test, y_pred_bernoulli_binary)
7 print(f'Binary Vectorizer with BernoulliNB Accuracy: {accuracy_bernoulli_binary}')
8 print(f'Binary Vectorizer with BernoulliNB Classification Report:\n{report_bernoulli_binary}')
```

Binary Vectorizer with BernoulliNB Accuracy: 0.7380952380952381

Binary Vectorizer with BernoulliNB Classification Report:

	precision	recall	f1-score	support
negative	0.81	0.62	0.70	21
positive	0.69	0.86	0.77	21
accuracy			0.74	42
macro avg	0.75	0.74	0.73	42
weighted avg	0.75	0.74	0.73	42

```

1 # Gaussian Naive Bayes with Binary Vectorizer
2 accuracy_gaussian_binary, report_gaussian_binary = apply_gaussian_nb(X_train_binary, X_test_binary, y_train, y_test)
3 print(f'Binary Vectorizer with GaussianNB Accuracy: {accuracy_gaussian_binary}')
4 print(f'Binary Vectorizer with GaussianNB Classification Report:\n{report_gaussian_binary}')
```

Binary Vectorizer with GaussianNB Accuracy: 0.7142857142857143

Binary Vectorizer with GaussianNB Classification Report:

	precision	recall	f1-score	support
negative	0.70	0.76	0.73	21
positive	0.74	0.67	0.70	21
accuracy			0.71	42
macro avg	0.72	0.71	0.71	42
weighted avg	0.72	0.71	0.71	42

```

1 # Save the Multinomial Naive Bayes model with TF-IDF
2 dump(nb_tfidf, 'nb_tfidf_model.joblib')

↩ ['nb_tfidf_model.joblib']

1 # Load the saved model
2 nb_tfidf = load('nb_tfidf_model.joblib')

1 # Predict a new sentence
2 def predict_sentence(model, vectorizer, sentence):
3     sentence_vectorized = vectorizer.transform([sentence])
4     prediction = model.predict(sentence_vectorized)
5     return prediction[0]

1 # Get input from the user
2 user_input = input("Enter a sentence to predict its class: ")
3
4 # Predict and display the result
5 prediction = predict_sentence(nb_tfidf, tfidf, user_input)
6 print(f"The predicted class for the sentence '{user_input}' is: {prediction}")
7
8 # Display the model's accuracy
9 print(f'The model accuracy is: {accuracy_tfidf}')

↩ Enter a sentence to predict its class: pressure in sports
The predicted class for the sentence 'pressure in sports' is: positive
The model accuracy is: 0.7619047619047619

```

Conclusion:

TF-IDF Vectorizer

- with MultinomialNB Accuracy: 0.7619047619047619
- with BernoulliNB Accuracy: 0.7380952380952381
- with GaussianNB Accuracy: 0.7380952380952381

Count Vectorizer

- with MultinomialNB Accuracy: 0.6904761904761905
- with BernoulliNB Accuracy: 0.7380952380952381
- with GaussianNB Accuracy: 0.7142857142857143

Binary Vectorizer

- with MultinomialNB Accuracy: 0.6904761904761905
- with BernoulliNB Accuracy: 0.7380952380952381
- with GaussianNB Accuracy: 0.7142857142857143

Highest Accuracy:

TF-IDF Vectorizer with MultinomialNB Accuracy: 0.7619047619047619