

> Importing the necessary Libraries

[ ] 1 cell hidden

Import Dataset

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

Choose Files

bbc-news-data.csv

- bbc-news-data.csv(text/csv) - 5080260 bytes, last modified: 7/31/2020 - 100% done

Saving bbc-news-data.csv to bbc-news-data.csv

```
1 import pandas as pd                #Imports the pandas library as pd
2 import numpy as np                 #Imports the numpy library for numerical computing
3
4 data_file="bbc-news-data.csv"
5
6 data = pd.read_csv(data_file,sep='\t')    ## Load the CSV file into a DataFrame
```

```
1 data.head(10)
```

Table icon

	category	filename	title	content	cleaned_content
0	business	001.txt	Ad sales boost Time Warner profit	Quarterly profits at US media giant TimeWarne...	quarterly profit u medium giant timewarner jum...
1	business	002.txt	Dollar gains on Greenspan speech	The dollar has hit its highest level against ...	dollar hit highest level euro almost three mon...
2	business	003.txt	Yukos unit buyer faces loan claim	The owners of embattled Russian oil giant Yuk...	owner embattled russian oil giant yukos ask bu...
3	business	004.txt	High fuel prices hit BA's profits	British Airways has blamed high fuel prices f...	british airway blamed high fuel price 40 drop ...
4	business	005.txt	Pernod takeover talk lifts Domecq	Shares in UK drinks and food firm Allied Dome...	share uk drink food firm allied domecq risen s...
5	business	006.txt	Japan narrowly escapes recession	Japan's economy teetered on the brink of a te...	japan economy teetered brink technical recessi...
6	business	007.txt	Jobs growth still slow in the US	The US created fewer jobs than expected in Ja...	u created fewer job expected january fall jobs...
7	business	008.txt	India calls for fair trade rules	India, which attends the G7 meeting of seven ...	india attends g7 meeting seven leading industr...
8	business	009.txt	Ethiopia's crop production up 24%	Ethiopia produced 14.27 million tonnes of cro...	ethiopia produced 14 27 million tonne crop 200...
9	business	010.txt	Court rejects \$280bn tobacco case	A US government claim accusing the country's ...	u government claim accusing country biggest to...

Grid icon

Bar chart icon

Next steps:

Generate code with data

View recommended plots

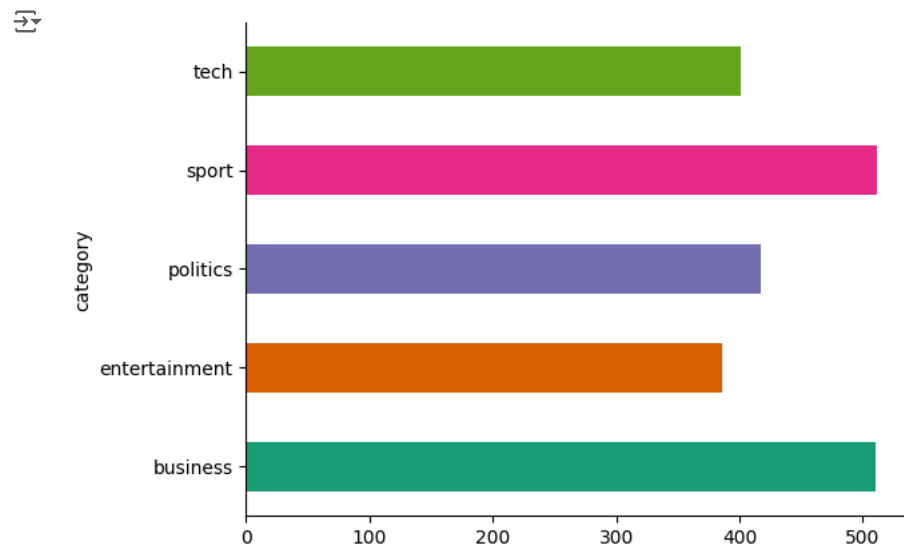
```
1 data.info()
```

Code icon

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2225 entries, 0 to 2224
Data columns (total 5 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   category         2225 non-null   object
1   filename         2225 non-null   object
2   title            2225 non-null   object
3   content          2225 non-null   object
4   cleaned_content  2225 non-null   object
dtypes: object(5)
memory usage: 87.0+ KB
```

category

```
1
2
3 # @title category
4
5 from matplotlib import pyplot as plt
6 import seaborn as sns
7 data.groupby('category').size().plot(kind='barh', color=sns.palettes.mpl_palette('Dark2'))
8 plt.gca().spines[['top', 'right']].set_visible(False)
```



## ✓ Text cleaning/Stop Words Removal

### Text Cleaning without stemming/stop word removal

```
1 #import re                                #Import module for regular expressions
2
3 # Text cleaning function
4 #def clean_text(text):
5 #     text = re.sub(r'\s+', ' ', text) # to remove extra spaces, newlines, and tabs.
6 #     text = re.sub(r'\W', ' ', text)  # to replace all non-word characters with a space.
7 #     return text.strip().lower()      # to remove any leading and trailing whitespace from the text and Converts the text to lowercase.
```

## ✓ Text Cleaning using stemming/stop word removal

```
1 # Text cleaning function
2 def clean_text(text):
3     text = re.sub(r'\W', ' ', text)
4     text = text.lower()
5     text = re.sub(r'\s+', ' ', text)
6     stop_words = set(stopwords.words('english'))
7     lemmatizer = WordNetLemmatizer()
8     #text = ' '.join([word for word in text.split() if word not in stop_words])
9     text = ' '.join([lemmatizer.lemmatize(word) for word in text.split() if word not in stop_words])
10    return text
11
12
13 # Clean the content
14 data['cleaned_content'] = data['content'].apply(clean_text)
```

## ✓ TFIDF Vectorization

### Max Features in Content(text) Column Unigram

```
1 # Extract the content column
2 content = data['content']
3
4 # Use CountVectorizer to count unique terms
5 vectorizer = CountVectorizer(stop_words='english')
6 X_counts = vectorizer.fit_transform(content)
7
8 # Get the number of unique terms
9 num_unique_terms = len(vectorizer.get_feature_names_out())
10
11 print(f'Number of unique terms (Unigrams): {num_unique_terms}')
```

➡ Number of unique terms (Unigrams): 28980

```
1 # Vectorization
2 tfidf_vectorizer = TfidfVectorizer(max_features=5000)
3 X = tfidf_vectorizer.fit_transform(data['cleaned_content'])
4 y = data['category']
```

## Splitting data into Train-test subset

```
1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state=42)
2
3
4 # Check the distribution of the categories in the test set
5 print("Category distribution in test set:")
6 print(y_test.value_counts())
```

Category distribution in test set:

category	
business	142
sport	122
tech	104
politics	95
entertainment	94
Name: count, dtype: int64	

## Model Training and Evaluation

```
1 from sklearn.metrics import accuracy_score, classification_report, precision_score, recall_score, confusion_matrix
2 import seaborn as sns
3
4 def train_and_evaluate_model(model, X_train, y_train, X_test, y_test):
5     model.fit(X_train, y_train)
6     y_pred = model.predict(X_test)
7     accuracy = accuracy_score(y_test, y_pred)
8     report = classification_report(y_test, y_pred, target_names=y_test.unique())
9     cm = confusion_matrix(y_test, y_pred, labels=y_test.unique())
10    return accuracy, report, cm, y_pred
11
12 # Calculate and print TP, FP, FN, TN
13 def calculate_metrics(cm):
14     for i in range(len(cm)):
15         tp = cm[i, i]
16         fp = cm[:, i].sum() - tp
17         fn = cm[i, :].sum() - tp
18         tn = cm.sum() - (tp + fp + fn)
19         print(f'Class {i}: TP: {tp}, FP: {fp}, TN: {tn}, FN: {fn}')
20
21 def plot_confusion_matrix(cm, classes):
22     plt.figure(figsize=(10, 7))
23     sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=classes, yticklabels=classes)
24     plt.ylabel('True label')
25     plt.xlabel('Predicted label')
26     plt.title('Confusion Matrix')
27     plt.show()
28
29 # Cross-validation
30 def cross_validate_model(model, X, y, cv=5):
31     scores = cross_val_score(model, X, y, cv=cv)
32     print(f"Cross-Validation Scores: {scores}")
33     print(f"Mean Accuracy: {scores.mean()}")
34     print(f"Standard Deviation: {scores.std()}")
```

## Naive Bayes

```
1 nb_model = MultinomialNB()
2 nb_accuracy, nb_report, nb_cm, nb_y_pred = train_and_evaluate_model(
3     nb_model, X_train, y_train, X_test, y_test)
```

### Evaluation\_nb

```
1 # Print results
2 print("Naive Bayes Classifier:")
3 print("Accuracy:", nb_accuracy)
4 print("Classification Report:\n", nb_report)
```

Naive Bayes Classifier:

Accuracy: 0.9658886894075404

Classification Report:

	precision	recall	f1-score	support
business	0.96	0.95	0.95	142
sport	0.99	0.95	0.97	94
politics	0.92	0.97	0.94	95
entertainment	1.00	0.99	1.00	122
tech	0.96	0.97	0.97	104
accuracy			0.97	557

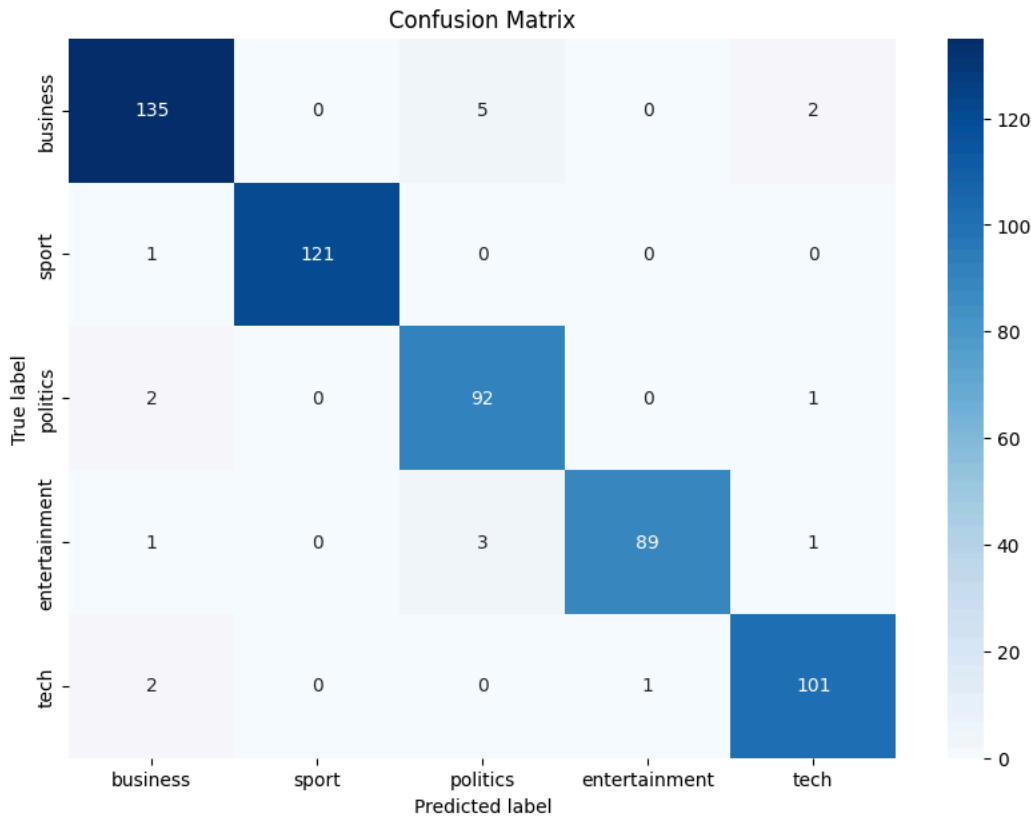
macro avg	0.97	0.97	0.97	557
weighted avg	0.97	0.97	0.97	557

```
1
2 print("\nMetrics per Class:")
3 calculate_metrics(nb_cm)
```



```
Metrics per Class:
Class 0: TP: 135, FP: 6, TN: 409, FN: 7
Class 1: TP: 121, FP: 0, TN: 435, FN: 1
Class 2: TP: 92, FP: 8, TN: 454, FN: 3
Class 3: TP: 89, FP: 1, TN: 462, FN: 5
Class 4: TP: 101, FP: 4, TN: 449, FN: 3
```

```
1 # Plot the confusion matrix
2 plot_confusion_matrix(nb_cm, classes=y_test.unique())
```



## Cross Validation\_nb

```
1 print("Naive Bayes Classifier (Cross-Validation):")
2 cross_validate_model(nb_model, X, y)
```



```
Naive Bayes Classifier (Cross-Validation):
Cross-Validation Scores: [0.96179775 0.95730337 0.93932584 0.98202247 0.98202247]
Mean Accuracy: 0.964494382022472
Standard Deviation: 0.016167286066161763
```

## Decision Tree

```
1 X_train_dense = X_train.toarray()
2 X_test_dense = X_test.toarray()

1 # Decision Tree
2 dt_model = DecisionTreeClassifier(random_state=42)
3 dt_accuracy, dt_report, dt_cm, dt_y_pred = train_and_evaluate_model(
4     dt_model, X_train, y_train, X_test, y_test)
```

## Evaluation\_dt

```
1 print("Decision Tree Classifier:")
2 print("Accuracy:", dt_accuracy)
3 print("Classification Report:\n", dt_report)
```



```
Decision Tree Classifier:
Accuracy: 0.8527827648114902
Classification Report:
```

precision recall f1-score support

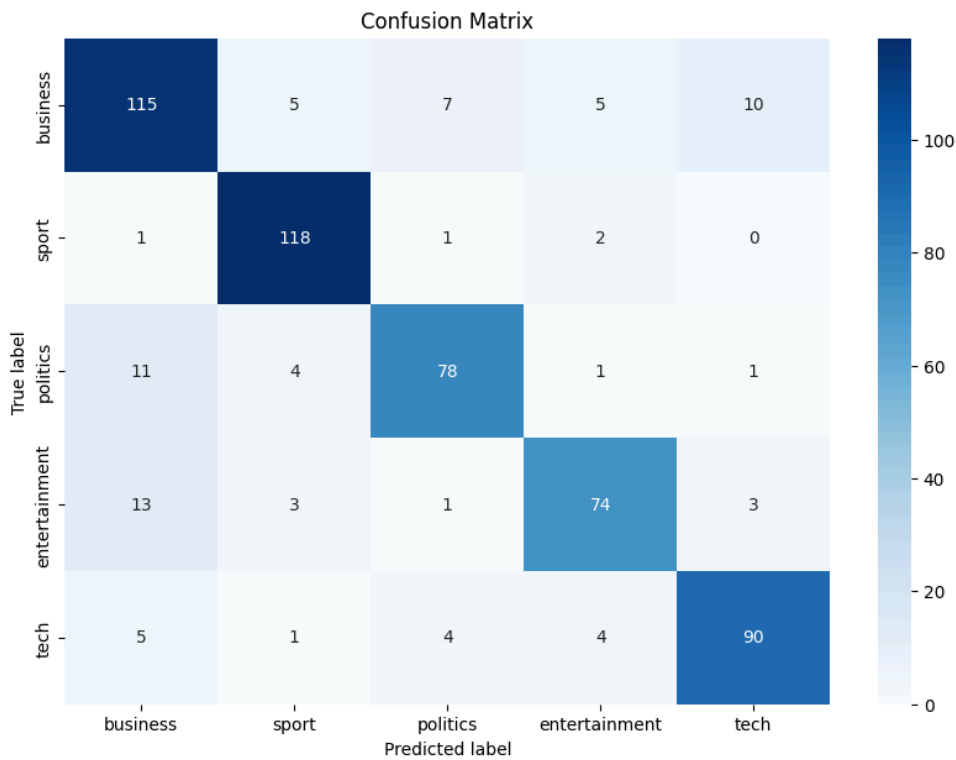
business	0.79	0.81	0.80	142
sport	0.86	0.79	0.82	94
politics	0.86	0.82	0.84	95
entertainment	0.90	0.97	0.93	122
tech	0.87	0.87	0.87	104
accuracy			0.85	557
macro avg	0.86	0.85	0.85	557
weighted avg	0.85	0.85	0.85	557

```
1 print("\nMetrics per Class:")
2 calculate_metrics(dt_cm)
```



```
Metrics per Class:
Class 0: TP: 115, FP: 30, TN: 385, FN: 27
Class 1: TP: 118, FP: 13, TN: 422, FN: 4
Class 2: TP: 78, FP: 13, TN: 449, FN: 17
Class 3: TP: 74, FP: 12, TN: 451, FN: 20
Class 4: TP: 90, FP: 14, TN: 439, FN: 14
```

```
1 # Plot the confusion matrix
2 plot_confusion_matrix(dt_cm, classes=y_test.unique())
```



## Cross Validation\_dt

```
1 print("Decision Tree Classifier (Cross-Validation):")
2 cross_validate_model(dt_model, X, y)
```



```
Decision Tree Classifier (Cross-Validation):
Cross-Validation Scores: [0.77078652 0.81123596 0.78876404 0.80449438 0.86741573]
Mean Accuracy: 0.8085393258426967
Standard Deviation: 0.032577294875362
```

## Random Forest

```
1 # Random Forest
2 rf_model = RandomForestClassifier(n_estimators=100, random_state=42)
3 rf_accuracy, rf_report, rf_cm, rf_y_pred = train_and_evaluate_model(
4     rf_model, X_train, y_train, X_test, y_test)
```

## Evaluation-rf

```
1 print("Random Forest Classifier:")
2 print("Accuracy:", rf_accuracy)
3 print("Classification Report:\n", rf_report)
```

Random Forest Classifier:

Accuracy: 0.9551166965888689

Classification Report:

	precision	recall	f1-score	support
business	0.94	0.95	0.95	142
sport	0.99	0.94	0.96	94
politics	0.93	0.94	0.93	95
entertainment	0.96	0.98	0.97	122
tech	0.96	0.96	0.96	104
accuracy			0.96	557
macro avg	0.96	0.95	0.95	557
weighted avg	0.96	0.96	0.96	557

```
1 print("\nMetrics per Class:")
2 calculate_metrics(rf_cm)
```

Metrics per Class:

Class 0: TP: 135, FP: 8, TN: 407, FN: 7

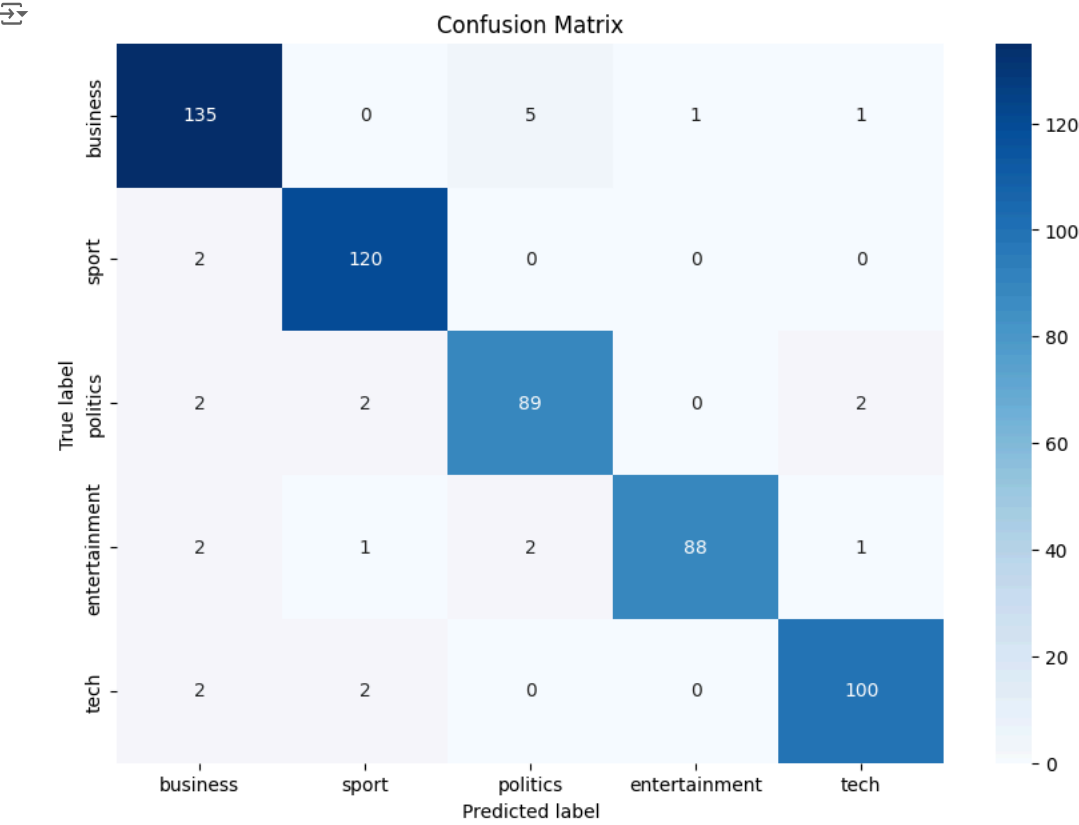
Class 1: TP: 120, FP: 5, TN: 430, FN: 2

Class 2: TP: 89, FP: 7, TN: 455, FN: 6

Class 3: TP: 88, FP: 1, TN: 462, FN: 6

Class 4: TP: 100, FP: 4, TN: 449, FN: 4

```
1 # Plot the confusion matrix
2 plot_confusion_matrix(rf_cm, classes=y_test.unique())
```



Cross Validation\_rf

```
1 print("Random Forest Classifier (Cross-Validation):")
2 cross_validate_model(rf_model, X, y)
3
```

Random Forest Classifier (Cross-Validation):

Cross-Validation Scores: [0.94606742 0.95955056 0.92359551 0.96629213 0.97303371]

Mean Accuracy: 0.9537078651685393


Standard Deviation: 0.01749925603050048

Comparison

```

1 # Store the results in a dictionary
2 results = {
3     "Model": ["Naive Bayes", "Decision Tree", "Random Forest"],
4     "Accuracy": [nb_accuracy, dt_accuracy, rf_accuracy],
5     "Classification Report": [nb_report, dt_report, rf_report]
6 }
7
8 # Convert the results dictionary to a DataFrame for better visualization
9 results_df = pd.DataFrame(results)
10 print("Comparison of Model Accuracies and Reports:")
11 print(results_df)
12
13
14 # Bar plot for accuracies
15 plt.figure(figsize=(8, 6))
16 plt.bar(results_df["Model"], results_df["Accuracy"], color=['blue', 'green', 'red'])
17 plt.xlabel("Model")
18 plt.ylabel("Accuracy")
19 plt.title("Accuracy Comparison of Different Models")
20 plt.show()
21

```



```

Comparison of Model Accuracies and Reports:

```

	Model	Accuracy		Classification Report		
0	Naive Bayes	0.965889	precision	recall	f1-score	...
1	Decision Tree	0.852783	precision	recall	f1-score	...
2	Random Forest	0.955117	precision	recall	f1-score	...

