

```
import pandas as pd
import matplotlib.pyplot as plt
import tensorflow as tf
import nltk
import sklearn
```

```
from tensorflow import keras
from keras.preprocessing.text import text_to_word_sequence
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
from nltk.stem import WordNetLemmatizer
from nltk.corpus import wordnet
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.model_selection import train_test_split
from sklearn import svm
from sklearn.naive_bayes import MultinomialNB
from sklearn import svm
from sklearn.naive_bayes import GaussianNB
from sklearn.metrics import roc_auc_score
from sklearn.pipeline import Pipeline
from sklearn.metrics import classification_report
from sklearn.metrics import accuracy_score
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
dataset = pd.read_csv('/content/BBC News.csv')
```

test-set = pd.read_csv('content/bbc_news.csv')

dataset.head()

	Article ID	Text	category
0	1833	worldcom ex-boss launches defence lawyer def...	Business
1	154	german business confidence slides german...	Business
2	1101	bbc poll indicates economic gloom citizen...	Business
3	1976	lifestyle governs mobile choice faster bett...	tech
4	917	enron bosses in \$168m payout eighteen former	Business.

target_category = dataset['category'].unique()

print(target_category)

['business' 'tech' 'politics' 'sport' 'entertainment']

dataset['category_id'] = dataset['category'].factorize()[0]

dataset.head()

	Article ID	Text	category	category_id
0	1833	worldcom ex-boss launches defence lawyers def...	Business	0
1	154	german business confidence slides german busin...	Business	0
2	1101	bbc poll indicates economic gloom citizens in...	Business	0
3	1976	lifestyle governs mobile choice-faster bett...	tech	1
4	917	enron bosses in \$168m payout eighteen former	Business	0

category = dataset[['category', 'category_id']].drop_duplicates(subset=['category_id'])

category

```

Category
0 business CategoryId
0
dataset.groupby('category').categoryId.count()

```

Category	Count
Business	336
entertainment	273
politics	274
Sport	346
tech	261

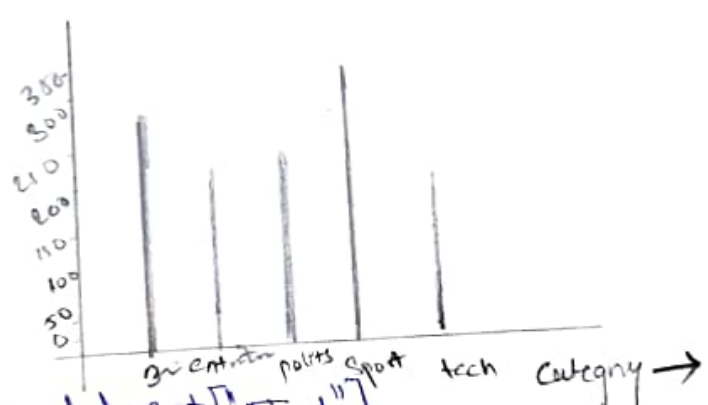
Name: categoryId dtype: int64

DATA VISUALIZATION

```

dataset.groupby('category').categoryId.count().plot.bar(ylim=0)

```



```

text = dataset["Text"]

```

```

text.head()

```

- 0 worldcom ex-boss launches defence lawyer def --
- 1 german business confidence slides. german busin..
- 2 bbc poll indicates economic gloom citizen in --
- 3 lifestyle governs mobile choice faster bett ---
- 4 enron bosses in \$168m payout eighteen formere ---

Name: Text dtype: object


```
print(classification_report(y_test_predict, y_test, target_names=target_category))
```

Naive Bayes Train Accuracy score: 98%

Naive Bayes Test Accuracy score: 96%

	precision	recall	f1-score	support
business	0.98	0.94	0.96	105
tech	0.89	1.00	0.94	73
politics	0.94	0.92	0.93	84
sport	1.00	0.98	0.99	106
entertainment	0.95	0.94	0.94	79
accuracy			0.96	447
macro avg	0.95	0.96	0.95	447
weighted avg	0.96	0.96	0.96	447

Decision Tree

```
dt = pipeline([['idf', TfidfVectorizer()], ('dt', DecisionTreeClassifier())])
```

```
dt.fit(x_train, y_train)
```

```
test_predict = dt.predict(x_test)
```

```
train_accuracy = round(dt.score(x_train, y_train)*100)
```

```
test_accuracy = round(accuracy_score(test_predict, y_test)*100)
```

```
print("Decision Tree Train Accuracy score: {}%".format(train_accuracy))
```

```
print("Decision Tree Test Accuracy score: {}%".format(test_accuracy))
```

```
print()
```

```
print(classification_report(test_predict, y_test, target_name=target_category))
```

Decision Tree Train Accuracy score: 100%

Decision Tree Test Accuracy score: 76%

	Precision	recall	f1-score	support
business	0.68	0.73	0.71	94
tech	0.68	0.74	0.71	76
politics	0.78	0.76	0.77	84
sport	0.91	0.79	0.84	121
entertainment	0.73	0.79	0.76	72
accuracy			0.76	447
macro avg	0.76	0.76	0.76	447
weighted avg	0.77	0.76	0.77	447

RANDOM FOREST CLASSIFIER

```
rfc = pipeline([('tfidf', TfidfVectorizer()), ('rfc', RandomForestClassifier(n_estimators=100))])
```

```
rfc.fit(x_train, y_train)
```

```
test_predict = rfc.predict(x_test)
```

```
train_accuracy = round(rfc.score(x_train, y_train)*100)
```

```
test_accuracy = round(accuracy_score(test_predict, y_test)*100)
```

```
print("K-Nearest Neighbour Train Accuracy score: {}%".format(train_accuracy))
```

```
print("K-Nearest Neighbour Test Accuracy score: {}%".format(test_accuracy))
```

```
print()
```

```
print(classification_report(test_predict, y_test, target_names=target_category))
```

```
K-Nearest Neighbour Train Accuracy score: 100%
```

```
K-Nearest Neighbour Test Accuracy score: 93%
```

	precision	recall	f1-score	support
business	0.97	0.88	0.92	112
tech	0.88	1.00	0.94	72
politics	0.85	0.93	0.89	75

Sport	1.00	0.93	0.96	112
entertainment	0.92	0.95	0.94	76
accuracy			0.93	447
macro avg	0.93	0.94	0.93	447
Weighted avg	0.94	0.93	0.93	447

TEST SET

test-set-head(1)

	Article ID	Text	category
0	1833	worldcom ex-boss launches defec	business
1	154	german business confidence	business
2	1101	bbc poll indicates	business
3	1976	lifestyle governs mobile choice	tech
4	914	emob bosses in \$100m payout	business

— * — * —