

```

1 import nltk
2 from nltk import word_tokenize
3 import sys
4
5 nltk.download('stopwords')
6 nltk.download('punkt')
7 nltk.download('omw-1.4')

[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package omw-1.4 to /root/nltk_data...
[nltk_data]   Package omw-1.4 is already up-to-date!
True

```

```

1 from google.colab import drive
2 import pandas as pd
3
4 drive.mount('/content/gdrive/', force_remount=True)

Mounted at /content/gdrive/

```

▼ Step 1

```

1 df = pd.read_csv('/content/gdrive/My Drive/Colab_Notebooks/federalist.csv')
2 df = df.astype({"author": 'category'})
3 print(df[:10])
4 authors = {}
5 for author in df['author']:
6     if author in authors.keys():
7         authors[author] = authors.get(author, 0) + 1
8     else:
9         authors[author] = 1
10 for author in authors:
11     print(author + " : " + str(authors.get(author)))

```

	author	text
0	HAMILTON	FEDERALIST. No. 1 General Introduction For the...
1	JAY	FEDERALIST No. 2 Concerning Dangers from Forei...
2	JAY	FEDERALIST No. 3 The Same Subject Continued (C...
3	JAY	FEDERALIST No. 4 The Same Subject Continued (C...
4	JAY	FEDERALIST No. 5 The Same Subject Continued (C...
5	HAMILTON	FEDERALIST No. 6 Concerning Dangers from Disse...
6	HAMILTON	FEDERALIST. No. 7 The Same Subject Continued (...
7	HAMILTON	FEDERALIST No. 8 The Consequences of Hostiliti...
8	HAMILTON	FEDERALIST No. 9 The Union as a Safeguard Agai...
9	MADISON	FEDERALIST No. 10 The Same Subject Continued (...

```
HAMILTON : 49
JAY : 5
MADISON : 15
HAMILTON AND MADISON : 3
HAMILTON OR MADISON : 11
```

▼ Step 2

```
1 from sklearn.model_selection import train_test_split
2 X = df.text.values
3 y = df.author.values
4 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, train_size=0.8)
5 print(X_train.shape)
6 print(X_test.shape)

(66,)
(17,)
```

▼ Step 3

```
1 from sklearn.feature_extraction.text import TfidfVectorizer
2 from nltk.corpus import stopwords
3 stopwords = set(stopwords.words('english'))
4 vectorizer = TfidfVectorizer(stop_words=stopwords)
5 X_train = vectorizer.fit_transform(X_train)
6 X_test = vectorizer.transform(X_test)
7 print(X_train.shape)
8 print(X_test.shape)

(66, 7876)
(17, 7876)
```

▼ Step 4

```
1 from sklearn.naive_bayes import BernoulliNB
2 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, cc
3
4 naive_bayes = BernoulliNB()
5 naive_bayes.fit(X_train, y_train)
6
7 pred = naive_bayes.predict(X_test)
8 print('accuracy score: ', accuracy_score(y_test, pred))
```

```
accuracy score: 0.5882352941176471
```

▼ Step 5

```
1 vectorizer_v2 = TfidfVectorizer(min_df=2, max_df=0.5, ngram_range=(1, 2), stop_words =
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, train_size=0.8
3 X_train = vectorizer_v2.fit_transform(X_train)
4 X_test = vectorizer_v2.transform(X_test)
5
6 naive_bayes.fit(X_train, y_train)
7
8 pred = naive_bayes.predict(X_test)
9 print('accuracy score: ', accuracy_score(y_test, pred))

accuracy score: 0.9411764705882353
```

▼ Step 6

▼ No Parameters

```
1 from sklearn.linear_model import LogisticRegression
2 classifier = LogisticRegression()
3 classifier.fit(X_train, y_train)
4
5 # evaluate
6 pred = classifier.predict(X_test)
7 print('accuracy score: ', accuracy_score(y_test, pred))

accuracy score: 0.5882352941176471
```

▼ With Parameters

```
1 classifier = LogisticRegression(class_weight='balanced', C = 2, )
2 classifier.fit(X_train, y_train)
3
4 # evaluate
5 pred = classifier.predict(X_test)
6 print('accuracy score: ', accuracy_score(y_test, pred))

accuracy score: 0.8235294117647058
```

▼ Step 7

```
1 from sklearn.neural_network import MLPClassifier
2 classifier = MLPClassifier(random_state = 1234, max_iter=400, hidden_layer_sizes=(30, 2
3 classifier.fit(X_train, y_train)
4 pred = classifier.predict(X_test)
5 print('accuracy score: ', accuracy_score(y_test, pred))
```

accuracy score: 0.8823529411764706

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

✓ 1s completed at 6:54 PM

