```
from google.colab import files
uploaded = files.upload()
     Choose Files | federalist.csv
    • federalist.csv(text/csv) - 1100616 bytes, last modified: 11/12/2022 - 100% done
     Saving federalist.csv to federalist.csv
import pandas as pd
df = pd.read_csv('federalist.csv', header=0, encoding='latin-1')
import nltk
nltk.download('stopwords')
     [nltk data] Downloading package stopwords to /root/nltk data...
                   Unzipping corpora/stopwords.zip.
     True
from nltk.corpus import stopwords
from sklearn.feature extraction.text import TfidfVectorizer
stopwords = set(stopwords.words('english'))
vectorizer = TfidfVectorizer(min df=2, max df=0.5, ngram range=(1,1))
X = df.text
y = df.author
X.head()
map = \{ \}
for i in df.author:
  if i in map:
    map[i] = map[i] + 1
  else:
    map[i] = 1
print(map)
     {'HAMILTON': 49, 'JAY': 5, 'MADISON': 15, 'HAMILTON AND MADISON': 3, 'HAMILTON O
from sklearn.model selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, train_size=0.
print(X train.shape)
```

```
print(X_test.shape)
    (66,)
     (17.)
# apply tfidf vectorizer
X_train = vectorizer.fit_transform(X_train) # fit and transform the train data
X_test = vectorizer.transform(X_test)
                                         # transform only the test data
# take a peek at the data
# this is a very sparse matrix because most of the 8613 words don't occur in each sms
print('train size:', X_train.shape)
print(X_train.toarray()[:5])
print('\ntest size:', X test.shape)
print(X_test.toarray()[:5])
    train size: (66, 4603)
    .0]]
                  0.03220702 0.
                                         ... 0.
                                                         0.
                                                                    0.
                                                                               1
     0.
                                         ... 0.
                  0.
                             0.
                                                         0.
                                                                    0.
                                                                               1
     [0.
                  0.
                             0.
                                         ... 0.
                                                                    0.
                                                         0.
                                                                               1
                                         ... 0.
     0.
                  0.
                             0.
                                                                    0.
     [0.
                  0.
                             0.
                                         ... 0.
                                                         0.05589859 0.
                                                                               11
    test size: (17, 4603)
    [[0.
                  0.
                             0.
                                         ... 0.
                                                         0.
                                                                    0.
     [0.
                  0.
                             0.
                                         ... 0.
                                                         0.02659749 0.
     [0.
                  0.
                             0.
                                         ... 0.
                                                         0.
                                                                    0.
                                                                               1
      [0.
                  0.
                             0.
                                         ... 0.
                                                         0.
                                                                    0.
                                                                               ]
     [0.
                  0.
                              0.
                                         ... 0.
                                                         0.
                                                                    0.
                                                                               ]]
from sklearn.naive bayes import MultinomialNB
from sklearn.naive bayes import BernoulliNB
naive bayes = BernoulliNB()
naive bayes.fit(X train, y train)
    BernoulliNB(alpha=1.0, binarize=0.0, class prior=None, fit prior=True)
prior p = sum(y train == 1)/len(y train)
from sklearn.metrics import accuracy score, precision score, recall score, f1 score, c
# make predictions on the test data
pred = naive bayes.predict(X test)
# print confusion matrix
```

```
print(confusion matrix(y test, pred))
    [[10 0 0 0]
     [3 0 0 0]
     [2 0 0 0]
     [2 0 0 0]]
print('accuracy score: ', accuracy_score(y_test, pred))
    accuracy score: 0.5882352941176471
vectorizer = TfidfVectorizer(min df=2, max df=0.5, max features=1000,ngram range=(1,2)
X = df.text
y = df.author
X.head()
    0
         FEDERALIST. No. 1 General Introduction For the...
         FEDERALIST No. 2 Concerning Dangers from Forei...
         FEDERALIST No. 3 The Same Subject Continued (C...
         FEDERALIST No. 4 The Same Subject Continued (C...
         FEDERALIST No. 5 The Same Subject Continued (C...
    Name: text, dtype: object
X train, X test, y train, y test = train test split(X, y, test size=0.2, train size=0.
print(X train.shape)
print(X test.shape)
    (66,)
    (17,)
# apply tfidf vectorizer
X train = vectorizer.fit transform(X train) # fit and transform the train data
                                            # transform only the test data
X test = vectorizer.transform(X test)
from sklearn.naive bayes import BernoulliNB
naive bayes = BernoulliNB()
naive bayes.fit(X train, y train)
    BernoulliNB(alpha=1.0, binarize=0.0, class prior=None, fit prior=True)
prior p = sum(y train == 1)/len(y train)
pred = naive bayes.predict(X test)
```

```
Author_Attribution_Amanuel_Shiferaw_aks210003 - Colaboratory
# print confusion matrix
print(confusion_matrix(y_test, pred))
     [[10 0
                0]
      [ 0 3
              0 01
      [ 1 0
             1 01
      [0 0 0 2]]
print('accuracy score: ', accuracy score(y test, pred))
     accuracy score: 0.9411764705882353
!pip install "scikit_learn==0.22.2.post1"
     Looking in indexes: <a href="https://pypi.org/simple">https://us-python.pkg.dev/colab-whee</a>
     Requirement already satisfied: scikit_learn==0.22.2.post1 in /usr/local/lib/pyth
    Requirement already satisfied: joblib>=0.11 in /usr/local/lib/python3.7/dist-pacl
```

import pandas as pd from sklearn.feature extraction.text import TfidfVectorizer from sklearn.model selection import train test split from sklearn.linear model import LogisticRegression from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score,] df = pd.read csv('federalist.csv', header=0, encoding='latin-1') # set up X and y X = df.textv = df.author# divide into train and test sets X train, X test, y train, y test = train test split(X, y, test size=0.2, train size=0. vectorizer = TfidfVectorizer(binary=True) X_train = vectorizer.fit_transform(X_train) # fit and transform the train data X test = vectorizer.transform(X test) # transform only the test data #train classifier = LogisticRegression(multi class='multinomial', solver='lbfqs', class weight classifier.fit(X train, y train) # evaluate pred = classifier.predict(X test)

Requirement already satisfied: numpy>=1.11.0 in /usr/local/lib/python3.7/dist-pac Requirement already satisfied: scipy>=0.17.0 in /usr/local/lib/python3.7/dist-pac

```
print('accuracy score: ', accuracy_score(y_test, pred))
    0. 50022520/1176/71
import pandas as pd
from sklearn.feature extraction.text import TfidfVectorizer
from sklearn.model_selection import train_test_split
from sklearn.linear model import LogisticRegression
from sklearn.metrics import accuracy score, precision score, recall score, f1 score, ]
from nltk.corpus import stopwords
from sklearn.feature extraction.text import TfidfVectorizer
df = pd.read_csv('federalist.csv', header=0, encoding='latin-1')
# set up X and y
X = df.text
y = df.author
from nltk.corpus import stopwords
from sklearn.feature_extraction.text import TfidfVectorizer
stopwords = set(stopwords.words('english'))
vectorizer = TfidfVectorizer(stop_words=stopwords, binary=True)
X = vectorizer.fit transform(df.text)
y = df.author
# divide into train and test
from sklearn.model selection import train test split
X train, X test, y train, y test = train test split(X, y, test size=0.2, train size=0.
from sklearn.neural network import MLPClassifier
classifier = MLPClassifier(solver='lbfgs', alpha=1e-5,
                   hidden_layer_sizes=(17, 2), random state=1)
classifier.fit(X train, y train)
from sklearn.metrics import accuracy score
from sklearn.metrics import precision score, recall score, f1 score
pred = classifier.predict(X test)
print('accuracy score: ', accuracy_score(y_test, pred))
    accuracy score: 0.5882352941176471
df = pd.read_csv('federalist.csv', header=0, encoding='latin-1')
# set up X and y
X = df.text
y = df.author
from nltk.corpus import stopwords
from sklearn.feature extraction.text import TfidfVectorizer
stopwords = set(stopwords.words('english'))
vectorizer = TfidfVectorizer(stop words=stopwords, binary=True)
```

```
X = vectorizer.fit transform(df.text)
y = df.author
# divide into train and test
from sklearn.model selection import train test split
X train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, train_size=0.
from sklearn.neural network import MLPClassifier
classifier = MLPClassifier(solver='lbfgs', alpha=1e-5,
                   hidden_layer_sizes=(10, 2), random state=1)
classifier.fit(X_train, y_train)
from sklearn.metrics import accuracy score
from sklearn.metrics import precision score, recall score, f1 score
pred = classifier.predict(X test)
print('accuracy score: ', accuracy_score(y_test, pred))
    accuracy score: 0.5882352941176471
df = pd.read_csv('federalist.csv', header=0, encoding='latin-1')
# set up X and y
X = df.text
y = df.author
from nltk.corpus import stopwords
from sklearn.feature extraction.text import TfidfVectorizer
stopwords = set(stopwords.words('english'))
vectorizer = TfidfVectorizer(stop words=stopwords, binary=True)
X = vectorizer.fit_transform(df.text)
y = df.author
# divide into train and test
from sklearn.model selection import train test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, train_size=0.
from sklearn.neural network import MLPClassifier
classifier = MLPClassifier(solver='lbfgs', alpha=1e-5,
                   hidden layer sizes=(20, 4), random state=1)
classifier.fit(X train, y train)
from sklearn.metrics import accuracy score
from sklearn.metrics import precision score, recall score, f1 score
pred = classifier.predict(X test)
print('accuracy score: ', accuracy score(y test, pred))
C→ accuracy score: 0.17647058823529413
```

Colab paid products - Cancel contracts here

✓ 0s completed at 10:28 PM

×