import pandas as pd import glob import seaborn as sns import matplotlib.pyplot as plt %matplotlib inline plt.style.use('ggplot') category list = ["sport", "world", "us", "business", "health", "entertainment", "sci tech"] directory\_list = ["data/sport/\*.txt", "data/world/\*.txt", "data/us/\*.txt", "data/business/\*.txt", "data/health/\*.txt", "data/entertainment/\*.txt", "data/sci tech/\*.txt",] text files = list(map(lambda x: glob.glob(x), directory list)) text files = [item for sublist in text files for item in sublist] len(text files) Out[2]: 1433 training data = [] for t in text files: f = open(t, 'r')f = f.read()t = f.split(' n')training data.append({'data' : t[0] + ' ' + t[1], 'flag' : category list.index(t[6]), 'category' : t[6]}) training data[0] Out[3]: {'data': "court agrees to expedite n.f.l.'s appeal the decision means a ruling could be made nearly two months before the regular season begins, time for the sides to work out a deal without delaying the season.", 'flag': 0, 'category': 'sport'} In [4]: training\_data[0:3] Out[4]: [{'data': "court agrees to expedite n.f.l.'s appeal the decision means a ruling could be made nearly two months before the regular season begins, time for the sides to work out a deal without delaying the season.", 'flag': 0, 'category': 'sport'}, {'data': 'no tsunami but fifa\'s corruption storm rages on though jack warner\'s threatened soccer "tsunami" r emains stuck in the doldrums, the corruption storm raging around fifa shows no sign of abating after another ex traordinary week for the game\'s governing body.', 'flag': 0, 'category': 'sport'}, {'data': 'chung backs bin hammam in fifa presidential race mohamed bin hammam received a boost to his campaign to become fifa president when he was backed by his former foe and ex-fifa executive committee member chung mong -joon tuesday.', 'flag': 0, 'category': 'sport'}] Create Data Frame training data = pd.DataFrame(training data, columns=['data', 'flag','category']) print(training data.shape) (1433, 3)training data.head() data flag category 0 court agrees to expedite n.f.l.'s appeal the d... sport 1 no tsunami but fifa's corruption storm rages o... sport 2 chung backs bin hammam in fifa presidential ra... 0 sport 3 rory mcilroy hangs on to slim lead at augusta ... 0 sport figure skating: south korean skater stumbles, ... 0 sport Count the frequency of each news articles in diff categories training\_data.flag.value\_counts() 352 280 256 3 2 194 5 167 6 118 66 Name: flag, dtype: int64 plt.figure(figsize=(8, 4), dpi=100, facecolor='w', edgecolor='k') sns.barplot(x=category\_list, y=training\_data.flag.value\_counts()) plt.show() 350 300 250 -200 150 100 50 0 world business health entertainment sci\_tech sport US Feature extraction from sklearn.model\_selection import train\_test\_split X\_train, X\_test, y\_train, y\_test = train\_test\_split(training\_data.data, training\_data.flag, random\_state=1) print(X\_train.shape, X\_test.shape, y\_train.shape, y\_test.shape) (1074,) (359,) (1074,) (359,)import pickle from sklearn.feature extraction.text import CountVectorizer #GET VECTOR COUNT  $\verb|count_vect| = CountVectorizer(strip_accents='ascii', token_pattern=u'(?ui) \verb|\b|\w*[a-z]+\w*\b', lowercase=True, lowercase$ X train cv = count vect.fit transform(X train) X test cv = count vect.transform(X test) **#SAVE WORD VECTOR** #pickle.dump(count\_vect.vocabulary\_, open("count\_vector.pkl","wb")) X train\_cv  $\text{Out}[11]: < 1074 \times 7156 \text{ sparse matrix of type '<class 'numpy.int64'>'}$ with 19847 stored elements in Compressed Sparse Row format> Fit the model and make predictions from sklearn.naive\_bayes import MultinomialNB naive\_bayes = MultinomialNB() naive\_bayes.fit(X\_train\_cv, y\_train) Out[12]: MultinomialNB() **Predictions** predictions\_nb = naive\_bayes.predict(X\_test\_cv) print(predictions\_nb)  $[5\ 1\ 6\ 1\ 3\ 3\ 2\ 6\ 5\ 0\ 3\ 1\ 3\ 3\ 3\ 1\ 0\ 0\ 0\ 0\ 0\ 5\ 3\ 6\ 0\ 0\ 5\ 6\ 0\ 3\ 2\ 0\ 2\ 0\ 0\ 5\ 1$  $5 \; 3 \; 3 \; 5 \; 0 \; 0 \; 1 \; 0 \; 0 \; 1 \; 5 \; 5 \; 1 \; 1 \; 3 \; 1 \; 0 \; 6 \; 0 \; 0 \; 1 \; 3 \; 3 \; 0 \; 2 \; 5 \; 1 \; 6 \; 0 \; 3 \; 1 \; 3 \; 0 \; 3 \; 3 \; 2 \; 5$  $1 \; 0 \; 0 \; 6 \; 0 \; 5 \; 5 \; 2 \; 0 \; 2 \; 0 \; 3 \; 2 \; 0 \; 1 \; 3 \; 1 \; 0 \; 0 \; 3 \; 3 \; 1 \; 3 \; 0 \; 0 \; 0 \; 5 \; 4 \; 4 \; 4 \; 2 \; 3 \; 6 \; 4 \; 4 \; 0 \; 3 \; 1$  $3\; 3\; 3\; 1\; 1\; 2\; 2\; 2\; 5\; 0\; 0\; 5\; 1\; 1\; 0\; 0\; 2\; 3\; 0\; 3\; 1\; 2\; 5\; 0\; 1\; 1\; 3\; 1\; 6\; 3\; 6\; 0\; 3\; 4\; 1\; 0\; 1$  $1 \; 2 \; 1 \; 0 \; 1 \; 1 \; 0 \; 0 \; 5 \; 1 \; 3 \; 2 \; 0 \; 3 \; 5 \; 1 \; 3 \; 3 \; 5 \; 3 \; 1 \; 1 \; 2 \; 5 \; 0 \; 0 \; 6 \; 3 \; 5 \; 2 \; 3 \; 5 \; 0 \; 0 \; 0 \; 6$ 1 3 1 2 0 3 3 3 1 1 1 2 2 3 4 0 2 0 5 2 1 3 2 4 0 5 0 3 0 0 5 1 6 3 0 0 0 0 2 0 3 0 5 0 6 1 0 0 1 0 3 1 0 5 1 2 1 5 2 0 1 2 0 0 1 0 3 0 3 0 2  $1 \; 2 \; 4 \; 3 \; 3 \; 0 \; 6 \; 6 \; 1 \; 0 \; 2 \; 0 \; 0 \; 1 \; 0 \; 1 \; 5 \; 1 \; 3 \; 1 \; 0 \; 6 \; 3 \; 2 \; 1 \; 0 \; 3 \; 2 \; 3 \; 3 \; 2 \; 1 \; 5 \; 1 \; 2 \; 0 \; 3$  $\begin{smallmatrix} 6 & 0 & 1 & 6 & 0 & 0 & 3 & 1 & 3 & 0 & 5 & 2 & 2 & 6 & 3 & 5 & 1 & 0 & 2 & 0 & 3 & 2 & 1 & 1 & 5 & 1 & 6 & 0 & 2 & 2 & 0 & 1 & 1 & 1 & 0 & 1 & 1 \\ \end{smallmatrix}$ 0 3 4 0 5 0 6 5 0 3 2 0 1 0 5 0 3 6 5 0 3 5 0 3 5 0] In [14]: df = pd.DataFrame({"Actual":y\_test, "Predicted":predictions\_nb}) df.head() Out[14]: **Actual Predicted** 1195 5 5 494 6 820 2 547 1035 3 3 category\_list Out[15]: ['sport', 'world', 'us', 'business', 'health', 'entertainment', 'sci\_tech'] df1 = df.replace(0,category\_list[0]) df2 = df1.replace(1,category\_list[1]) df3 = df2.replace(2,category\_list[2]) df4 = df3.replace(3,category\_list[3]) df5 = df4.replace(4,category\_list[4]) df6 = df5.replace(5,category\_list[5]) final\_df = df6.replace(6, category\_list[6]) final\_df.head(10) **Actual Predicted 1195** entertainment entertainment 494 world world 820 us sci\_tech 547 world world 1035 business business 48 sport business 824 us us 1335 sci\_tech sci\_tech 867 business entertainment 204 sport sport KNN Algorithm import pickle # make a k- nearest model from sklearn.neighbors import KNeighborsClassifier clf\_k = KNeighborsClassifier(n\_neighbors=32) model2 = clf\_k.fit(X\_train\_cv, y\_train) predictions = model2.predict(X\_test\_cv) print(predictions) #save model #pickle.dump(model2, open("model2.pkl", "wb")) Check the results for Naive-Bayes from sklearn.metrics import accuracy\_score, precision\_score, recall\_score,confusion\_matrix print("Accuracy score: ", accuracy\_score(y\_test, predictions\_nb)) print("Precision score: ", precision\_score(y\_test, predictions\_nb, pos\_label='positive', average='micro')) print("Recall score: ", recall score(y test, predictions nb, pos\_label='positive',average='micro')) Accuracy score: 0.7409470752089137 Precision score: 0.7409470752089137 Recall score: 0.7409470752089137 C:\Users\Asus\anaconda3\lib\site-packages\sklearn\metrics\\_classification.py:1295: UserWarning: Note that pos\_1 abel (set to 'positive') is ignored when average != 'binary' (got 'micro'). You may use labels=[pos\_label] to s pecify a single positive class. warnings.warn("Note that pos\_label (set to %r) is ignored when " C:\Users\Asus\anaconda3\lib\site-packages\sklearn\metrics\\_classification.py:1295: UserWarning: Note that pos\_1 abel (set to 'positive') is ignored when average != 'binary' (got 'micro'). You may use labels=[pos\_label] to s pecify a single positive class. warnings.warn("Note that pos\_label (set to %r) is ignored when " **Confusion Matrix** conf mat = confusion matrix( final\_df.Actual, final\_df.Predicted, labels=['sport', 'world', 'us', 'business', 'health', 'entertainmer conf\_mat Out[19]: array([[81, 2, 3, 3, 0, 1, 0], [2, 56, 1, 2, 0, 1, 0],[ 4, 11, 29, 3, 0, 5, 1], 1, 2, 9, 3, 4, 44, [ 4, 1, 3], [ 0, 1, 5, 2, 0], [ 9, 0, 1, 4, 0, 28, 0], 1, 0, 10, 0, 1, 19]], dtype=int64) [ 2, labels=['sport', 'world', 'us', 'business', 'health', 'entertainment', 'sci\_tech'] cm = confusion matrix(final df.Actual, final df.Predicted, labels) fig = plt.figure(figsize=(8, 6)) ax = fig.add\_subplot() sns.heatmap(cm, annot=True, fmt='g', ax=ax); plt.xlabel('Predicted') plt.ylabel('True') plt.show() C:\Users\Asus\anaconda3\lib\site-packages\sklearn\utils\validation.py:70: FutureWarning: Pass labels=['sport', 'world', 'us', 'business', 'health', 'entertainment', 'sci\_tech'] as keyword args. From version 1.0 (renaming o f 0.25) passing these as positional arguments will result in an error warnings.warn(f"Pass {args\_msg} as keyword args. From version " 81 - 70 60 50 4 1 2 40 30 5 9 0 2 - 20 0 0 28 0 - 10 19 0 10 0 Ó 5 Predicted Check the results for KNN from sklearn.metrics import accuracy score, precision score, recall score print("Accuracy score: ", accuracy score(y test, predictions)) print("Precision score: ", precision\_score(y\_test, predictions, pos\_label='positive',average='micro')) print("Recall score: ", recall score(y test, predictions, pos\_label='positive',average='micro')) Accuracy score: 0.17270194986072424 Precision score: 0.17270194986072424 Recall score: 0.17270194986072424 C:\Users\Asus\anaconda3\lib\site-packages\sklearn\metrics\ classification.py:1295: UserWarning: Note that pos 1 abel (set to 'positive') is ignored when average != 'binary' (got 'micro'). You may use labels=[pos\_label] to s pecify a single positive class. warnings.warn("Note that pos\_label (set to %r) is ignored when "

C:\Users\Asus\anaconda3\lib\site-packages\sklearn\metrics\\_classification.py:1295: UserWarning: Note that pos\_l abel (set to 'positive') is ignored when average != 'binary' (got 'micro'). You may use labels=[pos\_label] to s

pecify a single positive class.

\* HAPPY LEARNINIG \*\*

warnings.warn("Note that pos\_label (set to %r) is ignored when "