

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import os, os.path
4 import re
5 from collections import Counter
6 from random import seed
7 from random import randrange
8 from tabulate import tabulate
9 from sklearn.neighbors import KNeighborsClassifier
10 from sklearn.utils import shuffle
11 from sklearn.metrics import accuracy_score
12 from sklearn.naive_bayes import MultinomialNB
13 from sklearn import svm
14 from scipy.sparse import csr_matrix
15 from sklearn.model_selection import KFold
16 import nltk
17 nltk.download('stopwords')
18 from nltk.corpus import stopwords
19 from sklearn.preprocessing import MinMaxScaler
20 from sklearn.metrics import classification_report
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
```

```
1 from sklearn.feature_extraction.text import CountVectorizer
2 from sklearn.feature_extraction.text import TfidfTransformer
```

```
1 from google.colab import drive
2 drive.mount('/content/drive')
3 # !rm -r /content/drive/My\ Drive/Pattern\ Recognition/Project/dataset_large
4 # !apt-get install zip unzip
5 # !unzip /content/drive/My\ Drive/Pattern\ Recognition/Project/dataset_large.zip -d /conte
6 # ***** IF YOU WANT TO IMPORT DATASET FROM LOCAL MACHINE KINDLY PASS THE FOLDER LOCAT
7 datasetPath = "/content/drive/My Drive/Pattern Recognition/Project/dataset_large/"
```

[illegible]


```

12     x_train.append(data)
13     if (int(label) == 1):
14         X_test.append(1)
15     if (int(label) == 2):
16         X_test.append(1)
17     if (int(label) == 3):
18         X_test.append(1)
19     if (int(label) == 4):
20         X_test.append(1)
21     if (int(label) == 7):
22         X_test.append(2)
23     if (int(label) == 8):
24         X_test.append(2)
25     if (int(label) == 9):
26         X_test.append(2)
27     if (int(label) == 10):
28         X_test.append(2)
29     if files == "test":
30         #print(f)
31         y_train.append(data)
32         if (int(label) == 1):
33             y_test.append(1)
34         if (int(label) == 2):
35             y_test.append(1)
36         if (int(label) == 3):
37             y_test.append(1)
38         if (int(label) == 4):
39             y_test.append(1)
40         if (int(label) == 7):
41             y_test.append(2)
42         if (int(label) == 8):
43             y_test.append(2)
44         if (int(label) == 9):
45             y_test.append(2)
46         if (int(label) == 10):
47             y_test.append(2)
48
49
50

```

```

1 X_train = np.array(X_train)
2 X_test = np.array(y_test)
3 y_train = np.array(y_train)
4 y_test = np.array(y_test)
5 target_names = []
6 labels = np.unique(X_test)
7 for label in labels:
8     target_names.append('Rating '+str(label))
9 print(X_train.shape)
10 print(X_test.shape)
11 print(y_train.shape)

```

```
12 print(y_test.shape)
```

```
(25000,)  
(25000,)  
(25000,)  
(25000,)
```

```
1 import nltk  
2 nltk.download('wordnet')  
3 from nltk.stem import WordNetLemmatizer  
4  
5 def clean_dataset(X):  
6     documents = []  
7     stemmer = WordNetLemmatizer()  
8     for sen in range(0, len(X)):  
9  
10         # Remove all the special characters  
11         document = re.sub(r'\W', ' ', str(X[sen]))  
12  
13         # remove all single characters  
14         document = re.sub(r'\s+[a-zA-Z]\s+', ' ', document)  
15  
16         # Remove single characters from the start  
17         document = re.sub(r'^\s+[a-zA-Z]\s+', ' ', document)  
18  
19         # Substituting multiple spaces with single space  
20         document = re.sub(r'\s+', ' ', document, flags=re.I)  
21  
22         # Removing prefixed 'b'  
23         document = re.sub(r'^b\s+', '', document)  
24  
25         # Converting to Lowercase  
26         document = document.lower()  
27  
28         # Lemmatization  
29         document = document.split()  
30  
31         document = [stemmer.lemmatize(word) for word in document]  
32         document = ' '.join(document)  
33  
34         documents.append(document)  
35     return documents  
36 X_train = clean_dataset(X_train)  
37 y_train = clean_dataset(y_train)
```

```
[nltk_data] Downloading package wordnet to /root/nltk_data...  
[nltk_data]   Package wordnet is already up-to-date!
```

```
1 def transform(array):  
2     from sklearn.feature_extraction.text import CountVectorizer  
3     vectorizer = CountVectorizer(max_features=5000, min_df=5, max_df=0.7, stop_words=stopw
```

```

4     vectorized = vectorizer.fit_transform(array)
5     transformer = TfidfTransformer()
6     transformed = transformer.fit_transform(vectorized)
7     return np.array(transformed.toarray())
8

```

```

1 print("Creating Count Vectors for training data")
2 X_train = transform(X_train)
3 print(np.array(X_train).shape)
4 print("Creating Count Vectors for testing data")
5 y_train = transform(y_train)
6 print(np.array(y_train).shape)
7

```

```

    Creating Count Vectors for training data
    (25000, 5000)
    Creating Count Vectors for testing data
    (25000, 5000)

```

```

1 def pad_along_axis(array: np.ndarray, target_length: int, axis: int = 0):
2
3     pad_size = target_length - array.shape[axis]
4
5     if pad_size <= 0:
6         return array
7
8     npad = [(0, 0)] * array.ndim
9     npad[axis] = (0, pad_size)
10
11     return np.pad(array, pad_width=npad, mode='constant', constant_values=0)
12
13 X_train = pad_along_axis(X_train, y_train.shape[1], axis=1)

```

```

1 # convert to sparse matrix to dense matrix
2 print("Converting sparse matrix to dense matrix")
3 X_train_dense = csr_matrix(X_train)
4 y_train_dense = csr_matrix(y_train)

```

```

    Converting sparse matrix to dense matrix

```

```

1
2 print("Shuffling test and training datasets")
3 X_train_shuffled,X_test_shuffled = shuffle(X_train_dense,X_test)
4 print(X_train_shuffled.shape)
5 print(X_test_shuffled.shape)
6 y_train_shuffled,y_test_shuffled = shuffle(y_train_dense,y_test)
7 print(y_train_shuffled.shape)
8 print(y_test_shuffled.shape)
9
10

```

--

```
Shuffling test and training datasets
(25000, 5000)
(25000,)
(25000, 5000)
(25000,)
```

```
1 bestAverageK = 0
2 highestMean = 0
3 def mean(array,column,folds=5):
4     sum = 0
5     for i in range(1,folds+1):
6         sum = sum + float(array[i][column])
7     return round(sum / folds,3)
8
9 def formatMean(array):
10     global highestMean,bestAverageK
11     m = 0
12     formattedMean = ["Mean"]
13     for i in range(1,20):
14         m = mean(resultTable,i)
15
16         formattedMean.append(m)
17         if highestMean < m:
18             highestMean = m
19             bestAverageK = i
20     return formattedMean
21 resultTable = [[]]
22
23 kf = KFold(n_splits=5, random_state=None, shuffle=False)
24 for j in range(kf.get_n_splits(X_train_shuffled)):
25     for train_index, test_index in kf.split(X_train_shuffled):
26         split_input_train_dataset , split_label_train_dataset = X_train_shuffled[train_index]
27         split_input_test_dataset , split_label_test_dataset = X_train_shuffled[test_index],X
28
29     resultRow = [j+1]
30     for i in range(1 , 20):
31         knn = KNeighborsClassifier(n_neighbors=i)
32         knn.fit(split_input_train_dataset , split_label_train_dataset)
33         resultRow.append(str(round(accuracy_score(split_label_test_dataset,knn.predict(split
34     resultTable.append(resultRow)
35
36 #Calculate average for all value of K
37
38 resultTable.append(formatMean(resultTable))
39 table = tabulate(resultTable, headers=['Folds', "1","2","3","4","5","6","7","8","9","10","
40 print("Classification complete printing results")
41 print(table)
42 print(f"Highest mean is {highestMean} so the best K using 5 fold cross validation is {best
43 #bestAverageK =
44 #print(np.array(resultTable))
```

Classification complete printing results

Folds	1	2	3	4	5	6	7	8	9	10
1	0.526	0.527	0.535	0.525	0.539	0.54	0.54	0.531	0.542	0.537
2	0.526	0.527	0.535	0.525	0.539	0.54	0.54	0.531	0.542	0.537
3	0.526	0.527	0.535	0.525	0.539	0.54	0.54	0.531	0.542	0.537
4	0.526	0.527	0.535	0.525	0.539	0.54	0.54	0.531	0.542	0.537
5	0.526	0.527	0.535	0.525	0.539	0.54	0.54	0.531	0.542	0.537
Mean	0.526	0.527	0.535	0.525	0.539	0.54	0.54	0.531	0.542	0.537

Highest mean is 0.55 so the best K using 5 fold cross validation is 15



```

1 #using lowest k to classify test database
2 print(f"Classifying using KNN test dataset using best K {bestAverageK}")
3 clf = KNeighborsClassifier(n_neighbors=bestAverageK)
4 clf.fit(X_train_shuffled , X_test_shuffled)
5 knn_accuracy = accuracy_score(y_test_shuffled, clf.predict(y_train_shuffled))
6 print(f"Test dataset accuracy for best K {bestAverageK} is {knn_accuracy} ")
7 print(classification_report(y_test_shuffled, clf.predict(y_train_shuffled), target_names=t
8

```

Classifying using KNN test dataset using best K 15

Test dataset accuracy for best K 15 is 0.51252

	precision	recall	f1-score	support
Rating 1	0.51	0.49	0.50	12500
Rating 2	0.51	0.53	0.52	12500
accuracy			0.51	25000
macro avg	0.51	0.51	0.51	25000
weighted avg	0.51	0.51	0.51	25000

```

1 print(f"Classifying using multinomialNB")
2 clf = MultinomialNB()
3 clf.fit(X_train_shuffled , X_test_shuffled)
4 nb_accuracy = accuracy_score(y_test_shuffled, clf.predict(y_train_shuffled))
5 print(f"Test dataset accuracy is {nb_accuracy}")
6 print(classification_report(y_test_shuffled, clf.predict(y_train_shuffled), target_names=t
7

```

Classifying using multinomialNB

Test dataset accuracy is 0.5378

	precision	recall	f1-score	support
Rating 1	0.56	0.37	0.45	12500
Rating 2	0.53	0.70	0.60	12500
accuracy			0.54	25000
macro avg	0.54	0.54	0.52	25000
weighted avg	0.54	0.54	0.52	25000

```

1 print(f"Classifying using SVM")
2 clf = svm.SVC()
3 clf.fit(X_train_shuffled , X_test_shuffled)
4 svm_accuracy = accuracy_score(y_test_shuffled, clf.predict(y_train_shuffled))
5 print(f"Test dataset accuracy is {svm_accuracy} ")
6 print(classification_report(y_test_shuffled, clf.predict(y_train_shuffled), target_names=t
7

```

Classifying using SVM

Test dataset accuracy is 0.52848

	precision	recall	f1-score	support
Rating 1	0.53	0.53	0.53	12500
Rating 2	0.53	0.53	0.53	12500
accuracy			0.53	25000
macro avg	0.53	0.53	0.53	25000
weighted avg	0.53	0.53	0.53	25000