```
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 from sklearn.preprocessing import MinMaxScaler
17 from sklearn.metrics import classification_report
```

```
1 import pickle
2 import time
3
4 from sklearn.feature_extraction.text import TfidfVectorizer
5 from sklearn.decomposition import TruncatedSVD
6 from sklearn.pipeline import make_pipeline
7 from sklearn.preprocessing import Normalizer
8 from sklearn.neighbors import KNeighborsClassifier
9
```

```
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 LOC.
7 datasetPath = "/content/drive/My Drive/Pattern Recognition/Project/dataset_large/"
8 #datasetPath = "/content/drive/My Drive/Pattern Recognition/Project/dataset_large/"
```

```
    Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.moun
```

```
1 X_train = []
2 y_train = []
3 X_test = []
4 y_test = []
5 for files in os.listdir(datasetPath):
6   for f in os.listdir(datasetPath+files):
7     d = open(datasetPath+files+"/"+f, "r")
8     data = str(d.read())
```

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

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

```
10 print(y_train.shape)
11 print(X_test_raw.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'\^[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_raw = clean_dataset(X_train_raw)
37 X_test_raw = clean_dataset(X_test_raw)
```

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

```
 1 # Tfidf vectorizer:
```

```
 2 vectorizer = TfidfVectorizer(max_df=0.7, max_features=500,
 3                              min_df=5, stop_words='english',
 4                              use_idf=True)
 5
 6 # Build the tfidf vectorizer from the training data ("fit"), and apply it
 7 # ("transform").
 8 X_train_tfidf = vectorizer.fit_transform(X_train_raw)
 9
10 print("  Actual number of tfidf features: %d" % X_train_tfidf.get_shape()[1])
11
12 print("\nPerforming dimensionality reduction using LSA")
13 t0 = time.time()
14
15 # Project the tfidf vectors onto the first N principal components.
16 # Though this is significantly fewer features than the original tfidf vector,
17 # they are stronger features, and the accuracy is higher.
18 svd = TruncatedSVD(100)
19 lsa = make_pipeline(svd, Normalizer(copy=False))
20
21 # Run SVD on the training data, then project the training data.
22 X_train_lsa = lsa.fit_transform(X_train_tfidf)
23
24 print("  done in %.3fsec" % (time.time() - t0))
25
26 # explained_variance = svd.explained_variance_ratio_.sum()
27 # print("  Explained variance of the SVD step: {}%".format(int(explained_variance * 100)))
28
29
30 # Now apply the transformations to the test data as well.
31 X_test_tfidf = vectorizer.transform(X_test_raw)
32 X_test_lsa = lsa.transform(X_test_tfidf)
```

```
    Actual number of tfidf features: 500

  Performing dimensionality reduction using LSA
    done in 2.135sec
```

```
1 #scaling vector between 0 - 1
2 scaler = MinMaxScaler()
3 scaler.fit(X_train_lsa)
4 X_train_lsa_scaled =scaler.transform(X_train_lsa)
5
6 scaler = MinMaxScaler()
7 scaler.fit(X_test_lsa)
8 X_test_lsa_scaled =scaler.transform(X_test_lsa)
```

```
1 erageK = 0
2 tMean = 0
3 an(array,column,folds=5):
4 = 0
5 i in range(1,folds+1):
```

```
 6 n = sum + float(array[i][column])
 7 ·n round(sum / folds,3)
 8
 9 ·matMean(array):
10 al highestMean,bestAverageK
11 ·
12 attedMean = ["Mean"]
13 i in range(1,20):
14 = mean(resultTable,i)
15
16 ·mattedMean.append(m)
17  highestMean < m:
18 nighestMean = m
19 ·estAverageK = i
20 ·n formattedMean
21 Table = [[]]
22
23 ·old(n_splits=5, random_state=None, shuffle=False)
24 in range(kf.get_n_splits(X_train_tfidf)):
25 train_index, test_index in kf.split(X_train_tfidf):
26 split_input_train_dataset , split_label_train_dataset = X_train_tfidf[train_index],y_train|
27 split_input_test_dataset , split_label_test_dataset = X_train_tfidf[test_index],y_train[te:
28
29 ltRow = [j+1]
30 i in range(1 , 20):
31 <nn = KNeighborsClassifier(n_neighbors=i)
32 <nn.fit(split_input_train_dataset , split_label_train_dataset)
33 ·esultRow.append(str(round(accuracy_score(split_label_test_dataset,knn.predict(split_input_
34 ltTable.append(resultRow)
35
36 late average for all value of K
37
38 Table.append(formatMean(resultTable))
39 = tabulate(resultTable, headers=['Folds', "1","2","3","4","5","6","7","8","9","10","11","12
40 "Classification complete printing results")
41 table)
42 f"Highest mean is {highestMean} so the best K using 5 fold cross validation is {bestAverage
43
```

```
    Classification complete printing results
    | Folds   |  1  |  2  |  3   |  4   |  5   |  6  |  7  |  8   |  9   |  10
    |---------+-------+------+-------+-------+------+------+------+-------+-------+-------+
    |         |       |      |       |       |      |      |      |       |       |
    | 1       | 0.637 | 0.63 | 0.661 | 0.654 | 0.674 | 0.67 | 0.69 | 0.675 | 0.696 | 0.691
    | 2       | 0.637 | 0.63 | 0.661 | 0.654 | 0.674 | 0.67 | 0.69 | 0.675 | 0.696 | 0.691
    | 3       | 0.637 | 0.63 | 0.661 | 0.654 | 0.674 | 0.67 | 0.69 | 0.675 | 0.696 | 0.691
    | 4       | 0.637 | 0.63 | 0.661 | 0.654 | 0.674 | 0.67 | 0.69 | 0.675 | 0.696 | 0.691
    | 5       | 0.637 | 0.63 | 0.661 | 0.654 | 0.674 | 0.67 | 0.69 | 0.675 | 0.696 | 0.691
    | Mean    | 0.637 | 0.63 | 0.661 | 0.654 | 0.674 | 0.67 | 0.69 | 0.675 | 0.696 | 0.691
    Highest mean is 0.712 so the best K using 5 fold cross validation is 19
```

```
1 #Using best 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_lsa , y_train)
5 knn_accuracy = accuracy_score(y_test, clf.predict(X_test_lsa))
6 print(f"Test dataset accuracy for best K {bestAverageK} is {knn_accuracy} ")
7 print(classification_report(y_test, clf.predict(X_test_lsa), target_names=target_names))
8
```

```
    Classifying using KNN test dataset using best K 19
    Test dataset accuracy for best K 19 is 0.734
                  precision    recall  f1-score   support

        Rating 1       0.72      0.76      0.74     12500
        Rating 2       0.75      0.71      0.73     12500

        accuracy                           0.73     25000
       macro avg       0.73      0.73      0.73     25000
    weighted avg       0.73      0.73      0.73     25000
```

```
1 print(f"Classifying using multinomialNB")
2 clf = MultinomialNB()
3 clf.fit(X_train_lsa_scaled, y_train)
4 nb_accuracy = accuracy_score(y_test, clf.predict(X_test_lsa_scaled))
5 print(f"Test dataset accuracy is {nb_accuracy} ")
6 print(classification_report(y_test, clf.predict(X_test_lsa_scaled), target_names=target_na
7
```

```
    Classifying using multinomialNB
    Test dataset accuracy is 0.79676
                  precision    recall  f1-score   support

        Rating 1       0.82      0.76      0.79     12500
        Rating 2       0.78      0.83      0.80     12500

        accuracy                           0.80     25000
       macro avg       0.80      0.80      0.80     25000
    weighted avg       0.80      0.80      0.80     25000
```

```
1 print(f"Classifying using SVM")
2 clf = svm.SVC()
3 clf.fit(X_train_lsa_scaled , y_train)
4 svm_accuracy = accuracy_score(y_test, clf.predict(X_test_lsa_scaled))
5 print(f"Test dataset accuracy is {svm_accuracy} ")
6 print(classification_report(y_test, clf.predict(X_test_lsa_scaled), target_names=target_na
7
```

```
    Classifying using SVM
    Test dataset accuracy is 0.82336
                  precision    recall  f1-score   support
```

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| Rating 1     | 0.85      | 0.78   | 0.82     | 12500   |
| Rating 2     | 0.80      | 0.86   | 0.83     | 12500   |
|              |           |        |          |         |
| accuracy     |           |        | 0.82     | 25000   |
| macro avg    | 0.83      | 0.82   | 0.82     | 25000   |
| weighted avg | 0.83      | 0.82   | 0.82     | 25000   |