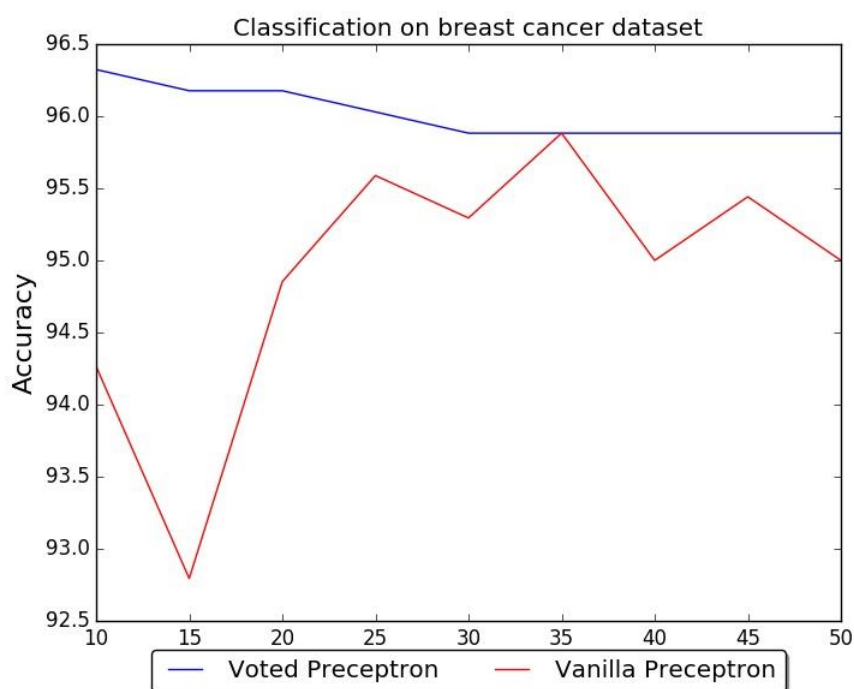


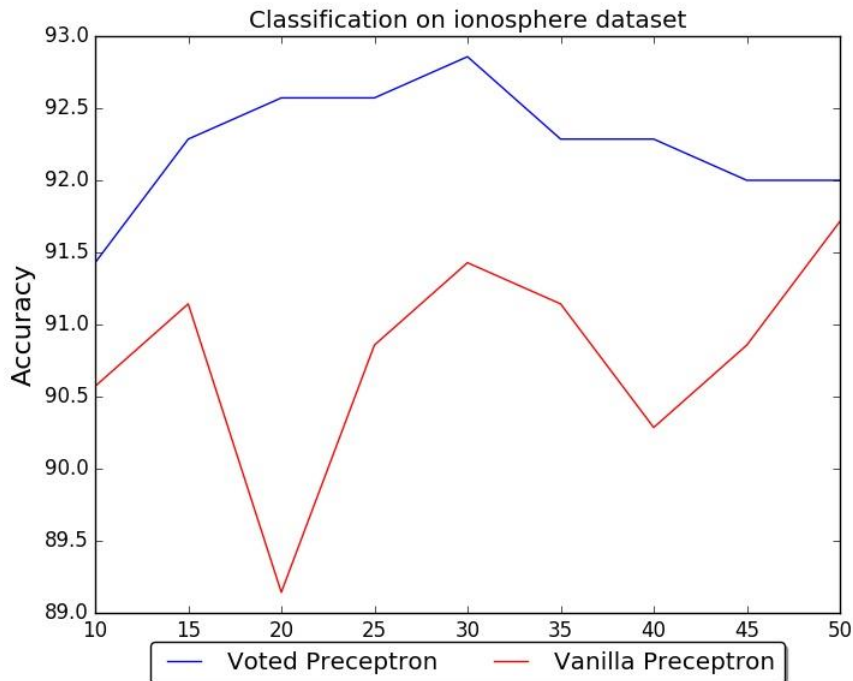
[Assignment-1]

Name: Ranajit Saha
Branch: M.Tech. CSE
Roll: 20172119

Q1: Compare the performance of Voted and Vanilla Perceptron Classification Models

Observation:





----Breast Cancer Dataset----

epoch = 10

Average Voted Perceptron Accuracy: 96.3235294118

Average Vanilla Perceptron Accuracy: 94.2647058824

epoch = 15

Average Voted Perceptron Accuracy: 96.1764705882

Average Vanilla Perceptron Accuracy: 92.7941176471

epoch = 20

Average Voted Perceptron Accuracy: 96.1764705882

Average Vanilla Perceptron Accuracy: 94.8529411765

epoch = 25

Average Voted Perceptron Accuracy: 96.0294117647

Average Vanilla Perceptron Accuracy: 95.5882352941

epoch = 30

Average Voted Perceptron Accuracy: 95.8823529412

Average Vanilla Perceptron Accuracy: 95.2941176471

epoch = 35

Average Voted Perceptron Accuracy: 95.8823529412

Average Vanilla Perceptron Accuracy: 95.8823529412

epoch = 40

Average Voted Perceptron Accuracy: 95.8823529412

Average Vanilla Perceptron Accuracy: 95.0

epoch = 45

Average Voted Perceptron Accuracy: 95.8823529412

Average Vanilla Perceptron Accuracy: 95.4411764706

epoch = 50

Average Voted Perceptron Accuracy: 95.8823529412

Average Vanilla Perceptron Accuracy: 95.0

----Ionosphere Dataset----

epoch = 10

Average Voted Perceptron Accuracy: 91.4285714286

Average Vanilla Perceptron Accuracy: 90.5714285714

epoch = 15

Average Voted Perceptron Accuracy: 92.2857142857

Average Vanilla Perceptron Accuracy: 91.1428571429

epoch = 20

Average Voted Perceptron Accuracy: 92.5714285714

Average Vanilla Perceptron Accuracy: 89.1428571429

epoch = 25

Average Voted Perceptron Accuracy: 92.5714285714

Average Vanilla Perceptron Accuracy: 90.8571428571

epoch = 30

Average Voted Perceptron Accuracy: 92.8571428571

Average Vanilla Perceptron Accuracy: 91.4285714286

epoch = 35

Average Voted Perceptron Accuracy: 92.2857142857

Average Vanilla Perceptron Accuracy: 91.1428571429

epoch = 40

Average Voted Perceptron Accuracy: 92.2857142857

Average Vanilla Perceptron Accuracy: 90.2857142857

epoch = 45

Average Voted Perceptron Accuracy: 92.0

Average Vanilla Perceptron Accuracy: 90.8571428571

epoch = 50

Average Voted Perceptron Accuracy: 92.0

Average Vanilla Perceptron Accuracy: 91.7142857143

Comments:

The accuracy of the prediction in case of voted perceptron is better than the vanilla perceptron. Also, the accuracy of the voted one is maintaining a value within very short range for different epochs, whereas the accuracy of vanilla varies a lot depending upon the epochs. We can see the accuracy of vanilla perceptron sometime dips too much, even after giving a very high accuracy for some epoch, if we increase the epoch further. But this is not the case with the voted perceptron, it is almost invariant with respect to the epochs after an epoch value. It seems the voted perceptron converges much earlier than vanilla.

Reason:

In perceptron training, we update the weights when the used weight fails to correctly classify the data point. In vanilla, whenever we find a misclassification, we update the hyperplane and train the model for other data points using the updated weight. But the problem is we don't keep track of how many examples has been correctly classified by the old hyperplane. So, when we move the hyperplane to a new position, it might be the case that the new hyperplane fails to correctly classify the points it was correctly classifying previously, and that number of misclassifications, say the set P , may be more than

the number of correctly classified examples. So, we lose the importance of the "more" correct classifier. When, we come back again for the next round of iteration, we once again have to adjust the hyperplane to make the set P to fall into the correct class. But, this way we cannot guarantee that the classifier accuracy will be improved with increased epoch, though it is guaranteed to converge at some point of time, obviously if the data is linearly classifiable.

On the other hand, the voted perceptron has the votes for each corresponding weight. This vote signifies which hyperplane classifies more examples correctly. So, when we predict the output we take the weighted value of the hyperplanes. This way the hyperplanes which were "more" correct and survives for a longer time has more contribution in predicting the output.

Code:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from random import shuffle
4 from sklearn.metrics import accuracy_score
5
6 def K_Fold_Cross_Validation(X, Y, K):
7     """
8     Divide the whole dataset into K-folds and
9     perform testing for each one the fold
10    while train the model using other K-1 folds
11    """
12    Y = Y.reshape(Y.shape[0],1)
13    sampleSize = int(X.shape[0] / K)
14    startPoint = 0
15    endPoint = startPoint + sampleSize
16    X_cv = X[startPoint:endPoint,:]
17    Y_cv = Y[startPoint:endPoint]
18    X_train = X[0:startPoint, :]
19    X_train = np.append(X_train, X[endPoint:,:],axis=0)
20    Y_train = Y[0:startPoint]
21    Y_train = np.append(Y_train, Y[endPoint:],axis=0)
22    divided_dataset = [[X_train, X_cv, Y_train, Y_cv]]
23    for i in range(K-1):
24        startPoint += sampleSize
25        endPoint=startPoint+sampleSize
26        X_cv = X[startPoint:endPoint,:]
27        Y_cv = Y[startPoint:endPoint]
28        X_train = X[0:startPoint,:]
29        X_train = np.append(X_train, X[endPoint:,:],axis=0)
30        Y_train = Y[0:startPoint]
31        Y_train = np.append(Y_train, Y[endPoint:],axis=0)
32        divided_dataset = np.append(divided_dataset, np.array([[X_train, X_cv, Y_train,Y_cv]]),axis=0)
33    return divided_dataset
34
35 def sign(val):
36     return -1 if val < 0 else 1
```

```

37
38 def votedPerceptron(X, Y, epoch):
39     """
40     Train the voted perceptron model
41     """
42     W = np.zeros((1, X.shape[1]))
43     B = np.zeros((1, 1))
44     C = np.zeros((1, 1))
45     n = 0
46     m = X.shape[0]
47     number_of_features = X.shape[1]
48     for round in range(epoch):
49         for i in range(m):
50             if Y[i] * (np.dot(W[n], X[i].reshape(number_of_features, 1)) + B[n]) <= 0:
51                 W = np.append(W, np.array([W[n] + Y[i] * X[i]]), axis=0)
52                 B = np.append(B, np.array([B[n] + Y[i]]), axis=0)
53                 C = np.append(C, [[1]], axis=0)
54                 n += 1
55             else:
56                 C[n] += 1
57         return W[1:, :], B[1:, :], C[1:, :]
58
59 def votedPerceptronPredict(X, Y, W, B, C):
60     """
61     Predict the output of the voted perceptron
62     using the voted weights found during training
63     """
64     Y_predicted = np.zeros((1, 1))
65     examples = Y.shape[0]
66     shape_w = W.shape
67     for e in range(examples):
68         temp_val = np.dot(W, X[e].T).reshape(shape_w[0], 1) + B
69         temp_val[temp_val > 0] = 1
70         temp_val[temp_val < 0] = -1
71         Y_predicted = np.append(Y_predicted, np.asarray([sign(np.dot(C.T, temp_val))]))
72     Y_predicted = Y_predicted[1:]
73     return Y_predicted.reshape(examples, 1)

```

```

74
75 def vanillaPerceptron(X, Y, epoch):
76     """
77     Train the vanilla perceptron model
78     """
79     bias = np.ones((X.shape[0], 1))
80     X = np.concatenate((X, bias), axis=1) # Append a column of 1s with X (for bias)
81     number_of_features = X[0].shape[0]
82     W = np.zeros((1, number_of_features))
83     m = X.shape[0]
84     for round in range(epoch):
85         for i in range(m):
86             if Y[i] * np.dot(W, X[i].reshape(number_of_features, 1)) <= 0:
87                 W = W + np.multiply(Y[i], X[i])
88     return W
89
90 def vanillaPerceptronPredict(X, Y, W):
91     """
92     Predict the output of the vanilla perceptron
93     using the weights found during training
94     """
95     bias = np.ones((X.shape[0], 1))
96     X = np.concatenate((X, bias), axis=1)
97     Y_predicted = np.dot(X, W.T)
98     Y_predicted[Y_predicted > 0] = 1
99     Y_predicted[Y_predicted < 0] = -1
100     return Y_predicted

```

```

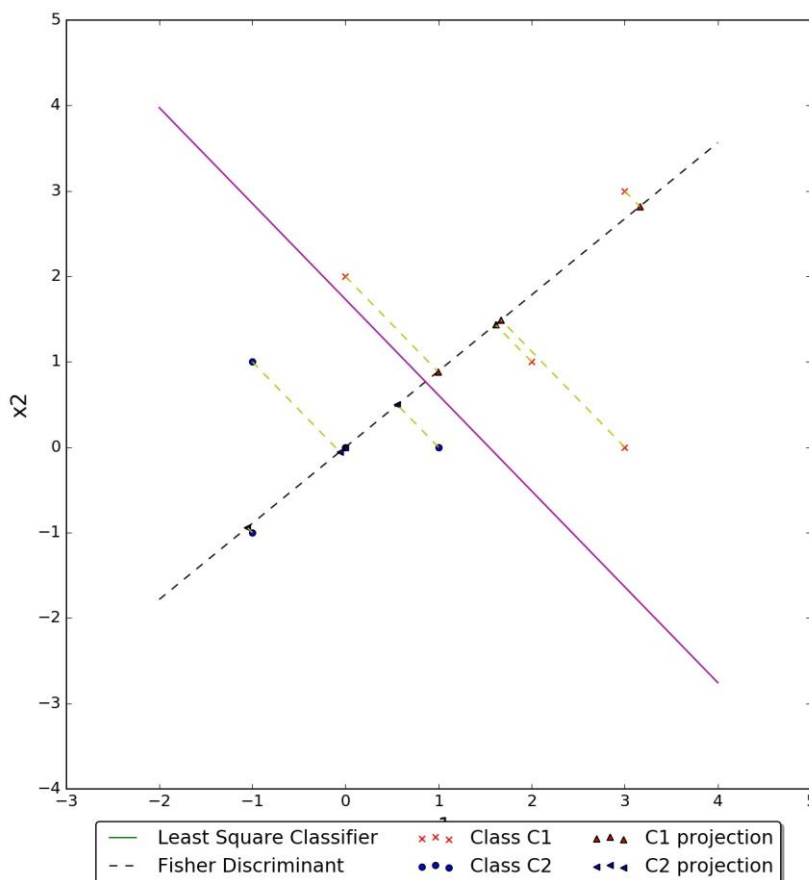
101
102 def perceptron(X, Y, epochs):
103     '''
104     Compare two perceptron algorithms
105     '''
106     K = 10 # For K-fold validation
107     accuracy_score_voted = []
108     accuracy_score_vanilla = []
109     divided_dataset = K_Fold_Cross_Validation(X,Y, K)
110     for epoch in epochs:
111         print("epoch = ", epoch)
112         accuracy_vanilla = 0
113         accuracy_voted = 0
114         for ij in range(divided_dataset.shape[0]):
115             X_train = divided_dataset[ij][0]
116             X_cv = divided_dataset[ij][1]
117             Y_train = divided_dataset[ij][2]
118             Y_cv = divided_dataset[ij][3]
119             W, B, C = votedPerceptron(X_train, Y_train, epoch)
120             Y_predicted_Voted = votedPerceptronPredict(X_cv, Y_cv, W, B, C)
121             accuracy_voted += accuracy_score(Y_cv, Y_predicted_Voted, normalize=True)
122             W = vanillaPerceptron(X_train, Y_train, epoch)
123             Y_predicted_Vanilla = vanillaPerceptronPredict(X_cv, Y_cv, W)
124             accuracy_vanilla += accuracy_score(Y_cv, Y_predicted_Vanilla, normalize=True)
125         print("Average Voted Perceptron Accuracy: ", accuracy_voted*100/K)
126         print("Average Vanilla Perceptron Accuracy: ", accuracy_vanilla*100/K)
127         accuracy_score_voted.append(accuracy_voted*100/K)
128         accuracy_score_vanilla.append(accuracy_vanilla*100/K)
129     return accuracy_score_voted, accuracy_score_vanilla
130

```

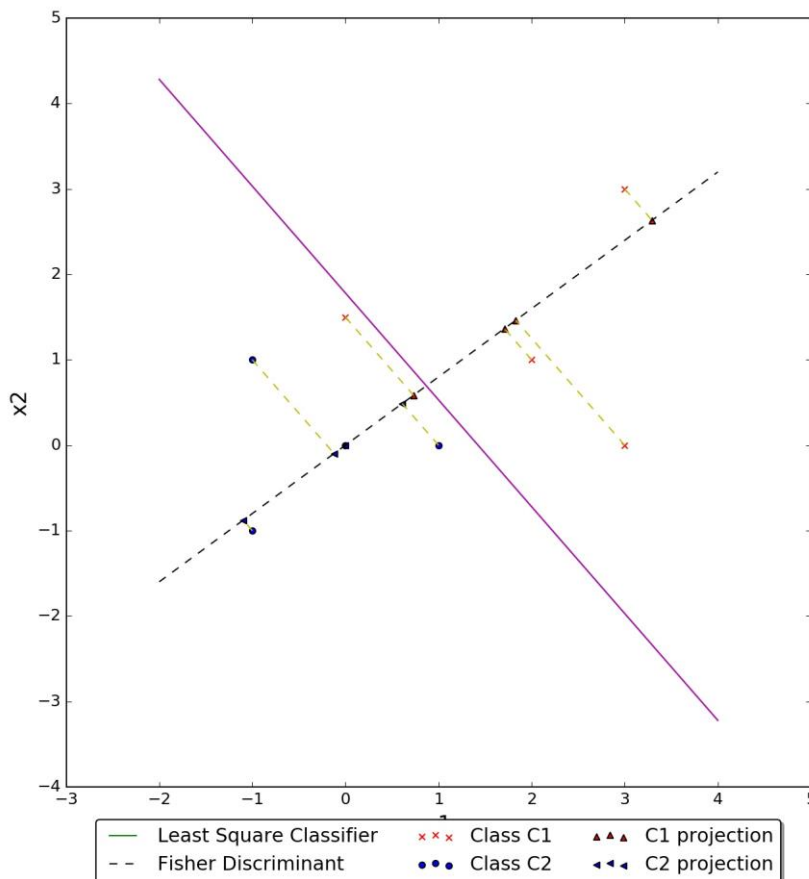
Q2. Comparison of classifiers - Least Square Method vs Fisher's LDA

Observation:

Fisher's discriminant vs Least square - C1 vs C2 - Table 1



Fisher's discriminant vs Least square - C1 vs C2 - Table 2



Comment:

The classifiers learnt using least square method and Fisher's LDA are the same. Direction of discriminant gives the direction of projection of data points in which the separation between the classes is maximum.

Reasons:

Fisher's LDA gives the direction of projection on which the data is reduced to a lower dimension and those projected data points are efficiently separable. Then we can apply any classifier on the projected points to discriminate them. The motive of Fisher's LDA is to project the data in such a way that they are also classifiable in the projected lower dimensional space and that would give the same decision boundary. This is observed in the output.

Code:

```

1 import numpy as np
2 import csv
3 import matplotlib.pyplot as plt
4
5 def LSC_binary_classification(X, Y):
6     """
7     Find the weight using
8     least square method
9     as  $W = (X.T * X)^{-1} * X.T * Y$ 
10    """
11    bias = np.ones((X.shape[0],1))
12    X_augmented = np.concatenate((bias, X),axis=1)
13    W = np.dot(np.linalg.pinv(np.dot(X_augmented.T, X_augmented)), np.dot(X_augmented.T, Y))
14    return W
15
16 def FisherLDA_binary_classification(X, Y):
17    """
18    Find the direction of projection which
19    maximises the class separation, the projected points
20    and the weights of the classifier
21    Direction of projection =  $S_W^{-1} * (mean1 - mean0)$ 
22    """
23    dataset = np.concatenate((X, Y), axis=1)
24    classes = np.unique(Y)
25    class_mean = []
26    num_of_dim = 1
27    for c in classes:
28        class_mean.append(((1 / len(Y[Y==c])) * np.sum(dataset[dataset[:, -1]==c], axis = 0))[: -1])
29    X_class1 = dataset[dataset[:, -1]==classes[0]]
30    X_class2 = dataset[dataset[:, -1]==classes[1]]
31    zero_mean_X_1 = X_class1[:, :-1] - class_mean[0]
32    zero_mean_X_2 = X_class2[:, :-1] - class_mean[1]
33    # Using Fisher's LDA formula
34    within_class_scatter = np.dot(zero_mean_X_1.T, zero_mean_X_1) + np.dot(zero_mean_X_2.T, zero_mean_X_2)
35    W = np.dot(np.linalg.inv(within_class_scatter), (class_mean[1] - class_mean[0]))
36    W = W / np.linalg.norm(W) # Unit vector along the direction of projection
37    projected_data_scalar = np.dot(X, W) # Magnitude of the projected points
38    projections = np.multiply(W.reshape(1,W.shape[0]),
39    projected_data_scalar.reshape(projected_data_scalar.shape[0], 1)) # Magnitude * Unit vector gives the actual vector
40    weight = LSC_binary_classification(projections, Y) # Classify the projected 1-D points using least square method
41    return projections, W, weight
42

```

```

43 def classify_LSC_LDA(X, Y, dataset_name):
44     """
45     Compare two classifiers using
46     1. Least square method
47     2. Fisher LDA
48     """
49     x_line = np.linspace(-2, 4)
50     plt.figure(figsize=(10,10))
51     plt.suptitle("Fisher's discriminant vs Least square - C1 vs C2 - " + dataset_name, fontsize = 18)
52     plt.xlabel("x1", fontsize = 16)
53     plt.ylabel("x2", fontsize = 16)
54     plt.scatter(X[:,4, 0], X[:,4, 1], c='r', marker='x', label='Class C1')
55     plt.scatter(X[:,4, 0], X[:,4, 1], c='b', marker='o', label='Class C2')
56     # Plot Least Square Method
57     LSC_weights = LSC_binary_classification(X, Y)
58     boundary_LSC = -LSC_weights[0] / LSC_weights[2] - (LSC_weights[1] / LSC_weights[2]) * x_line
59     plt.plot(x_line, boundary_LSC, color='g', label='Least Square Classifier')
60     # Plot LDA
61     projections, direction, LDA_weights = FisherLDA_binary_classification(X, Y)
62     direction_line = (direction[1] / direction[0]) * x_line
63     boundary_LDA = -LDA_weights[0] / LDA_weights[2] - (LDA_weights[1] / LDA_weights[2]) * x_line
64     # Plot Projections
65     plt.scatter(projections[:,4, 0], projections[:,4, 1], c='r', marker='^', label = "C1 projection")
66     plt.scatter(projections[:,4, 0], projections[:,4, 1], c='b', marker='<', label = "C2 projection")
67     for i in range(len(projections)):
68         plt.plot([X[i][0], projections[i][0]], [X[i][1], projections[i][1]], 'y--')
69     # Plot Discriminant and Boundary
70     plt.plot(x_line, direction_line, '--', color='black', label = "Fisher Discriminant")
71     plt.plot(x_line, boundary_LDA, color='m', label = "Fisher LDA")
72     plt.legend(loc='upper center', bbox_to_anchor=(0.5, -0.03), fancybox=True, shadow=True, ncol=3)
73
74     plt.savefig(dataset_name + '.jpg', dpi = 150)
75     # plt.show()
76     plt.clf()
77
78 #Table 1
79 X = np.array([[3, 3], [3, 0], [2, 1], [0, 2], [-1, 1], [0, 0], [-1, -1], [1, 0]])
80 Y = np.array([1, 1], [1], [1], [1], [-1], [-1], [-1], [-1]))
81 classify_LSC_LDA(X, Y, 'Table 1')
82
83 #Table 2
84 X = np.array([[3, 3], [3, 0], [2, 1], [0, 1.5], [-1, 1], [0, 0], [-1, -1], [1, 0]])
85 Y = np.array([1, 1], [1], [1], [1], [-1], [-1], [-1], [-1]))
86 classify_LSC_LDA(X, Y, 'Table 2')

```

Q3. Latent Semantic Analysis

- Classify documents using one-vs-rest voted perceptron after applying PCA:

Observation:

The accuracy of this method when we take 1110 singular values: 93.91634980988593

****Note:** Value of k depends on how many train documents we take

When we vary the value of K, the following accuracies have been reported:

```
raja@Ranajit-Inspiron-5559:~/Raja/Sem2/SMAI/Assignment1/FinalSubmission$ python3 q3_a.py '../LSAdata (copy)/train/' '../LSAdata (copy)/test/'
Accuracy using k = 10
82.50950570342205
Accuracy using k = 60
82.12927756653993
Accuracy using k = 110
82.12927756653993
Accuracy using k = 160
86.31178707224335
Accuracy using k = 210
88.212927756654
Accuracy using k = 260
91.25475285171103
Accuracy using k = 310
90.11406844106465
Accuracy using k = 360
89.35361216730038
Accuracy using k = 410
89.35361216730038
Accuracy using k = 460
91.25475285171103
Accuracy using k = 510
90.8745247148289
Accuracy using k = 560
90.8745247148289
Accuracy using k = 610
89.73384030418251
Accuracy using k = 660
90.11406844106465
Accuracy using k = 710
90.49429657794677
Accuracy using k = 760
92.01520912547528
Accuracy using k = 810
91.63498098859316
Accuracy using k = 860
90.11406844106465
Accuracy using k = 910
92.39543726235742
Accuracy using k = 960
93.5361216730038
Accuracy using k = 1010
92.39543726235742
Accuracy using k = 1060
92.01520912547528
Accuracy using k = 1110
93.91634980988593
Accuracy using k = 1160
93.5361216730038
Accuracy using k = 1210
93.15589353612167
Accuracy using k = 1260
93.15589353612167
Accuracy using k = 1310
93.15589353612167
Accuracy using k = 1360
93.5361216730038
Accuracy using k = 1410
92.77566539923954
Accuracy using k = 1460
92.39543726235742
raja@Ranajit-Inspiron-5559:~/Raja/Sem2/SMAI/Assignment1/FinalSubmission$
```

Code:

```

16 def sign(val):
17     return -1 if val < 0 else 1
18
19 def formTFMatrix(documents, terms):
20     '''
21     Construct the bag of words matrix for each document
22     where each row represents each document and
23     columns represent each terms
24     '''
25     tf_matrix = []
26     for words in documents:
27         doc_terms = []
28         for t in terms:
29             doc_terms.append(words.count(t))
30         tf_matrix.append(doc_terms)
31     return np.array(tf_matrix)
32
33 def multiclass_voted_perceptron(X, Y, epoch):
34     '''
35     Voted perceptron training, trained with for class
36     '''
37     m = X.shape[0]
38     number_of_features = X.shape[1]
39     W = np.zeros((1, number_of_features))
40     B = np.zeros((1, 1))
41     C = np.zeros((1, 1))
42     n = 0
43     for round in range(epoch):
44         for i in range(m):
45             if Y[i] * (np.dot(W[n], X[i].reshape(number_of_features, 1)) + B[n]) <= 0:
46                 W = np.append(W, np.array(W[n] + Y[i] * X[i]), axis=0)
47                 B = np.append(B, np.array([B[n] + Y[i]]), axis=0)
48                 C = np.append(C, [[1]], axis=0)
49                 n += 1
50             else:
51                 C[n] += 1
52     return W[1:, :], B[1:, :], C[1:, :]

```

```

53
54 def votedPerceptronPredict(X, Y, W_dict, B_dict, C_dict, class_labels):
55     """
56     Predict the class label using the trained voted perceptron model
57     Predict the class for which the absolute value prediction is maximum of other class
58     """
59     Y_predicted = np.zeros((1, 1))
60     examples = Y.shape[0]
61     for e in range(examples):
62         predict = np.asarray([])
63         for x in class_labels:
64             W = W_dict[x]
65             B = B_dict[x]
66             C = C_dict[x]
67             shape_w = W.shape
68             temp_val = np.dot(W, X[e].T).reshape(shape_w[0], 1) + B
69             # temp_val[temp_val>0] = 1
70             # temp_val[temp_val<0] = -1
71             predict = np.append(predict, np.dot(C.T, temp_val))
72     Y_predicted = np.append(Y_predicted, np.asarray([np.argmax(predict)]))
73
74     Y_predicted = Y_predicted[1:]
75     return Y_predicted.reshape(examples, 1)
76
77 def one_vs_all_voted(X, Y, X_test, Y_test, epoch):
78     """
79     Call voted perceptron for each of the classes
80     """
81     class_labels = np.unique(Y)
82     weight = {}
83     bias = {}
84     vote = {}
85     Y_original = deepcopy(Y)
86     for x in range(class_labels.shape[0]):
87         Y = deepcopy(Y_original)
88         for i in range(Y.shape[0]):
89             if Y[i][0] == x:
90                 Y[i][0] = 1
91             else:
92                 Y[i][0] = -1
93         weight[x], bias[x], vote[x] = multiclass_voted_perceptron(X, Y, epoch)
94     Y_predicted = votedPerceptronPredict(X_test, Y_test, weight, bias, vote, class_labels)
95     correct_class = 0
96     for i in range(Y_predicted.shape[0]):
97         if Y_predicted[i][0] == Y_test[i][0]:
98             correct_class += 1
99
100     for i in range(Y.shape[0]):
101         if Y[i][0] == x:
102             Y[i][0] = 1
103         else:
104             Y[i][0] = -1
105     weight[x], bias[x], vote[x] = multiclass_voted_perceptron(X, Y, epoch)
106     Y_predicted = votedPerceptronPredict(X_test, Y_test, weight, bias, vote, class_labels)
107     correct_class = 0
108     for i in range(Y_predicted.shape[0]):
109         if Y_predicted[i][0] == Y_test[i][0]:
110             correct_class += 1
111     print((correct_class / Y_predicted.shape[0]) * 100)
112
113 def extractDocs(dataPath):
114     """
115     Extract words from the documents
116     """
117     documents = []
118     class_labels = []
119     stop = set(stopwords.words('english')) # Removes the stop words eg the, an, a, and etc.
120     stemmer = PorterStemmer() # Converts the words to their root words
121     tokenizer = RegexpTokenizer(r'\w+') # Tokenize the document to get rid off non alphabets
122     folders = [dataPath + name + '/' for name in os.listdir(dataPath) if os.path.isdir(dataPath + name)]
123     class_name = [int(name) for name in os.listdir(dataPath) if os.path.isdir(dataPath + name)]
124     it = 0
125     for folder in folders:
126         docFiles = [folder + docName for docName in os.listdir(folder)]
127         doc_words = []
128         for f in docFiles:
129             file = codecs.open(f, 'r', 'ISO-8859-1')
130             file_content = []
131             for document in file:
132                 word_list = [stemmer.stem(line) for line in tokenizer.tokenize(document.lower()) if line not in '']
133                 file_content.extend([w for w in word_list if w not in stop])
134             file.close()
135             documents.append(file_content)
136             class_labels.append(class_name[it])
137             it += 1
138     return documents, (np.array([class_labels])).T

```

```

127
128 def getUniqueWords(documents):
129     '''
130     Find the dictionary of words
131     '''
132     words = [word for doc in documents for word in doc]
133     unique_words = list(set(words))
134     return unique_words
135
136 def classify_document_by_perceptron(dataPath, testPath):
137     '''
138     Classifies the documents after applying PCA
139     and then using the voted perceptron model
140     '''
141     avg_accuracy = []
142     documents, classes = extractDocs(dataPath)
143     test_documents, test_classes = extractDocs(testPath)
144     train_data_length = len(documents)
145     for test_doc in test_documents:
146         documents.append(test_doc)
147     terms = getUniqueWords(documents)
148     tf_matrix = formTFMatrix(documents, terms)
149     tf = TfidfTransformer(norm='l2', use_idf=True, smooth_idf=True, sublinear_tf=False)
150     tf_idf_matrix = tf.fit_transform(tf_matrix).todense() # Form the tf-idf matrix from bag of words
151     test_data = tf_idf_matrix[train_data_length:, :]
152     tf_idf_matrix = tf_idf_matrix[:train_data_length, :]
153     U, Sigma_temp, V_T = np.linalg.svd(tf_idf_matrix) # Singular Value Decomposition
154     Sigma = np.zeros((Sigma_temp.shape[0], Sigma_temp.shape[0]))
155     for i in range(Sigma_temp.shape[0]):
156         Sigma[i][i] = Sigma_temp[i]
157     # k list = list(range(10, min(tf_idf_matrix.shape[0], tf_idf_matrix.shape[1]), 50))
158     k_list = [1110]
159     for k in k_list:
160         correct_similarity = 0
161         # Reduce the dimension using the eigen value which contributes the most
162         reduced_data = np.dot(np.dot(U[:, :k], Sigma[:, :k]), V_T[:, :k])
163         epoch = 40 # Epoch for the perceptron training
164         print("Accuracy using k = ", k)
165         one_vs_all_voted(reduced_data, classes, test_data, test_classes, epoch)
166
167 dataPath = sys.argv[1] # './LSAdata (copy)/train/'
168 testPath = sys.argv[2] # './LSAdata (copy)/test/'
169
170 classify_document_by_perceptron(dataPath, testPath)
171

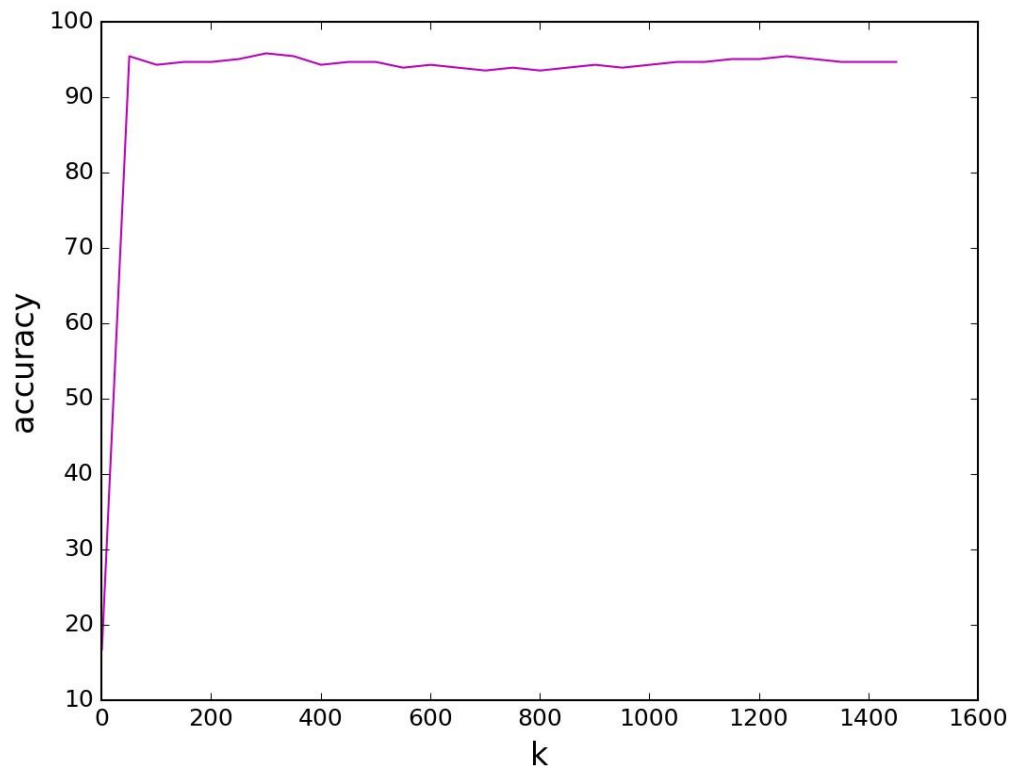
```

- Predict the class label of the test document using cosine similarity:

Observation:

As we increase the value of k, after some value the accuracy vs the values of k graph saturates, but at k = 1251 the accuracy is maximum.

Dimension (k) vs Accuracy



SVD Accuracy:

Accuracy for k = 1

16.730038022813687

Accuracy for k = 51

95.43726235741445

Accuracy for k = 101

94.29657794676807

Accuracy for k = 151

94.67680608365019

Accuracy for k = 201

94.67680608365019

Accuracy for k = 251

95.05703422053232

Accuracy for k = 301

95.81749049429658

Accuracy for k = 351

95.43726235741445

Accuracy for k = 401

94.29657794676807

Accuracy for k = 451

94.67680608365019

Accuracy for k = 501

94.67680608365019

Accuracy for k = 551

93.91634980988593

Accuracy for k = 601

94.29657794676807

Accuracy for k = 651
93.91634980988593
Accuracy for k = 701
93.5361216730038
Accuracy for k = 751
93.91634980988593
Accuracy for k = 801
93.5361216730038
Accuracy for k = 851
93.91634980988593
Accuracy for k = 901
94.29657794676807
Accuracy for k = 951
93.91634980988593
Accuracy for k = 1001
94.29657794676807
Accuracy for k = 1051
94.67680608365019
Accuracy for k = 1101
94.67680608365019
Accuracy for k = 1151
95.05703422053232
Accuracy for k = 1201
95.05703422053232
Accuracy for k = 1251
95.43726235741445
Accuracy for k = 1301
95.05703422053232
Accuracy for k = 1351
94.67680608365019
Accuracy for k = 1401
94.67680608365019
Accuracy for k = 1451
94.67680608365019

Code:

```

15
16 def extractDocs(dataPath):
17     """
18     Extract words from the documents
19     """
20     documents = []
21     class_labels = []
22     stop = set(stopwords.words('english')) # Removes the stop words eg the, an, a, and etc.
23     stemmer = PorterStemmer() # Converts the words to their root words
24     tokenizer = RegexpTokenizer(r'\w+') # Tokenize the document to get rid of non alphabets
25     folders = [dataPath + name + '/' for name in os.listdir(dataPath) if os.path.isdir(dataPath + name)]
26     class_name = [int(name) for name in os.listdir(dataPath) if os.path.isdir(dataPath + name)]
27     it = 0
28     for folder in folders:
29         docFiles = [folder + docName for docName in os.listdir(folder)]
30         doc_words = []
31         for f in docFiles:
32             file = codecs.open(f, 'r', 'ISO-8859-1')
33             file_content = []
34             for document in file:
35                 word_list = [stemmer.stem(line) for line in tokenizer.tokenize(document.lower()) if line not in '']
36                 file_content.extend([w for w in word_list if w not in stop])
37             file.close()
38             documents.append(file_content)
39             class_labels.append(class_name[it])
40             it += 1
41     return documents, (np.array([class_labels])).T
42
43 def getUniqueWords(documents):
44     """
45     Find the dictionary of words
46     """
47     words = [word for doc in documents for word in doc]
48     unique_words = list(set(words))
49     return unique_words
50
51 def formTFMatrix(documents, terms):
52     """
53     Construct the bag of words matrix for each document
54     where each row represents each document and
55     columns represent each terms
56     """
57     td_matrix = []
58     for words in documents:
59         doc_terms = []
60         for t in terms:

```

```

50
51 def formTFMatrix(documents, terms):
52     """
53     Construct the bag of words matrix for each document
54     where each row represents each document and
55     columns represent each terms
56     """
57     td_matrix = []
58     for words in documents:
59         doc_terms = []
60         for t in terms:
61             doc_terms.append(words.count(t))
62         td_matrix.append(doc_terms)
63     return np.array(td_matrix)
64
65 def readtestDoc(testPath):
66     """
67     Reading the test document
68     """
69     documents = []
70     stop = set(stopwords.words('english'))
71     stemmer = PorterStemmer()
72     tokenizer = RegexpTokenizer(r'\w+')
73     file = codecs.open(testPath, 'r', 'ISO-8859-1')
74     file_content = []
75     for document in file:
76         word_list = [stemmer.stem(line) for line in tokenizer.tokenize(document.lower()) if line not in '']
77         file_content.extend([w for w in word_list if w not in stop])
78     file.close()
79     documents.append(file_content)
80     return documents
81
82 def measure_similarity_by_cosine(dataPath, testPath, testDocLabel):
83     avg_accuracy = []
84     documents, classes = extractDocs(dataPath)
85     # test documents, test_classes = extractDocs(testPath)
86     test_documents = readtestDoc(testPath)
87     train_data_length = len(documents)
88     for test_doc in test_documents:
89         documents.append(test_doc)
90     terms = getUniqueWords(documents)
91     tf_matrix = formTFMatrix(documents, terms)
92     tf = TfidfTransformer(norm='l2', use_idf=True, smooth_idf=True, sublinear_tf=False)
93     tf_idf_matrix = tf.fit_transform(tf_matrix).todense() # Form the tf-idf matrix from bag of words
94     test_data = tf_idf_matrix[train_data_length:]

```

```

95 tf_idf_matrix = tf_idf_matrix[:train_data_length, :]
96 U, Sigma_temp, V_T = np.linalg.svd(tf_idf_matrix) # Singular Value Decomposition
97 Sigma = np.zeros((Sigma_temp.shape[0], Sigma_temp.shape[0]))
98 for i in range(Sigma_temp.shape[0]):
99     Sigma[i][i] = Sigma_temp[i]
100
101 # k list = list(range(1, min(tf_idf_matrix.shape[0], tf_idf_matrix.shape[1]), 50))
102 k_list = [1251] # Best value of k
103 # print("SVD Accuracy:")
104 for k in k_list:
105     correct_similarity = 0
106     reduced_data = np.dot(np.dot(U[:, :k], Sigma[:, :k]), V_T[:, :k])
107     cosine_similarity_matrix = sklearn.metrics.pairwise.cosine_similarity(reduced_data, test_data)
108     for x in range(test_data.shape[0]):
109         augmented_similarity_matrix = np.concatenate(
110             (cosine_similarity_matrix[:, x].reshape(cosine_similarity_matrix[:, 0].shape[0], 1),
111              classes.reshape(classes.shape[0], 1)), axis = 1)
112         sorted_similarity_matrix = augmented_similarity_matrix[augmented_similarity_matrix[:, 0].argsort()]
113         sorted_similarity_matrix = sorted_similarity_matrix[:, :-1]
114         predicted_label = -1
115         max_match = -1
116         top_10_results = sorted_similarity_matrix[:10, -1]
117         labels = np.unique(top_10_results)
118         for l in labels:
119             if max_match < list(top_10_results).count(l):
120                 max_match = list(top_10_results).count(l)
121                 predicted_label = l
122         print("The predicted class label of the document: ", predicted_label)
123         # if int(predicted_label) == int(test_classes[x]):
124             # correct_similarity += 1
125         # avg_accuracy.append(correct_similarity * 100 / test_data.shape[0])
126         # print("Accuracy for k = ", k),
127         # print(correct_similarity * 100 / test_data.shape[0])
128     # return k_list, avg_accuracy
129
130 dataPath = sys.argv[1] # './LSAdata (copy)/train/'
131 testPath = sys.argv[2] # './LSAdata (copy)/test/'
132 testDocLabel = sys.argv[3]
133 measure_similarity_by_cosine(dataPath, testPath, testDocLabel)
134 # k, accuracies = measure_similarity_by_cosine(dataPath, testPath, testDocLabel)
135 # plt.suptitle("Dimension (k) vs Accuracy", fontsize = 18)
136 # plt.xlabel("k", fontsize = 16)
137 # plt.ylabel("accuracy", fontsize = 16)
138 # plt.plot(k, accuracies, color='m', label = "Accracy for different k")
139 # plt.savefig('KvsAccuracy compressed.jpg', dpi = 150)

```