

1. For $k = 1, 5, 9$ and 15 , build k -nearest neighbor classifiers from the training data. For each of these values of k , write down a table of training errors (error on the training data) and the validation errors (error on the validation data). Which of these classifiers performs the best on validation data? What is the test error of this classifier?

	$k = 1$	$k = 5$	$k = 9$	$k = 15$
Training Error	0.0	0.0565	0.0685	0.093
Validation Error	0.082	0.101	0.102	0.104

The running time for calculating each training error is around 7 minutes and the running time for calculating each validation error is around 3 minutes.

According to this table, the classifier that performs the best on validation data is the 1-nearest neighbor classifier and the test error of this classifier is 0.094.

2. Project the training, validation and test data onto the column space of this matrix, and repeat part (1) of the problem. For $k = 1, 5, 9, 15$, write down a table of the training and validation errors, as well as the test error of the classifier which performs best on the validation data. How is the classification accuracy affected by projection? How does the running time of your program change when you run it on projected data?

	$k = 1$	$k = 5$	$k = 9$	$k = 15$
Training Error	0.0	0.192	0.23	0.259
Validation Error	0.32	0.299	0.296	0.295

The running time for calculating each training error is around 40 seconds and the running time for calculating each validation error is around 20 seconds.

According to this table, the classifier that performs the best on validation data is the 15-nearest neighbor classifier and the test error of this classifier is 0.303.

Comparing the data from two tables, we can see that matrix projection significantly reduces the running time of the program at the cost of accuracy.

Also, for the second part, although 15-nearest neighbor classifier performs best on the validation data, 9-nearest neighbor classifiers actually gives me better test errors, which further proves the fact that matrix projection reduces the accuracy.

```
1  import argparse
2  import numpy
3  import random
4
5  class Solver():
6      def __init__(self):
7          # read the arguments
8          parser = argparse.ArgumentParser()
9          parser.add_argument("training_set")
10         parser.add_argument("validation_set")
11         parser.add_argument("test_set")
12         parser.add_argument("projection_matrix")
13         args = parser.parse_args()
14         self.training_file = args.training_set
15         self.validation_file = args.validation_set
16         self.test_file = args.test_set
17         self.projection_file = args.projection_matrix
18         # list that contains the training data
19         self.training_data = []
20         # list that contains the validation data
21         self.validation_data = []
22         # list that contains the test data
23         self.test_data = []
24         # list that contains the rows of the projection matrix
25         self.projection_data = []
26     def load(self):
27         # open the training data file
28         with open(self.training_file) as training_file:
29             for line in training_file:
30                 data = line.split()
31                 for i in range(0, 785):
32                     data[i] = int(data[i])
33                 # add the data to the training data list
34                 self.training_data.append(data)
35         print("There are ", len(self.training_data), " training data")
36         # open the validation data file
37         with open(self.validation_file) as validation_file:
38             for line in validation_file:
39                 data = line.split()
40                 for i in range(0, 785):
41                     data[i] = int(data[i])
```

```
42         # add the data to the validation data list
43         self.validation_data.append(data)
44     print("There are ", len(self.validation_data), " validation data")
45     # open the test data file
46     with open(self.test_file) as test_file:
47         for line in test_file:
48             data = line.split()
49             for i in range(0, 785):
50                 data[i] = int(data[i])
51             # add the data to the test data list
52             self.test_data.append(data)
53     print("There are ", len(self.test_data), " test data")
54     # open the projection matrix file
55     with open(self.projection_file) as projection_file:
56         for line in projection_file:
57             data = line.split()
58             for i in range(0, 20):
59                 data[i] = float(data[i])
60             # add the data to the matrix data list
61             self.projection_data.append(data)
62     print("There are ", len(self.projection_data), " rows in the projection matrix")
63 def project(self):
64     m0 = numpy.array(self.projection_data)
65     # project the training data
66     training = []
67     for i in self.training_data:
68         training.append(i[:-1])
69     m1 = numpy.array(training)
70     m2 = numpy.matmul(m1, m0)
71     training_proj = []
72     for i in range(0, 2000):
73         r = m2[i].tolist()
74         r.append(self.training_data[i] [-1])
75         training_proj.append(r)
76     self.training_data = training_proj
77     # project the validation data
78     validation = []
79     for i in self.validation_data:
80         validation.append(i[:-1])
81     m1 = numpy.array(validation)
```

```
82     m2 = numpy.matmul(m1, m0)
83     validation_proj = []
84     for i in range(0, 1000):
85         r = m2[i].tolist()
86         r.append(self.validation_data[i][-1])
87         validation_proj.append(r)
88     self.validation_data = validation_proj
89     # project the test data
90     test = []
91     for i in self.test_data:
92         test.append(i[:-1])
93     m1 = numpy.array(test)
94     m2 = numpy.matmul(m1, m0)
95     test_proj = []
96     for i in range(0, 1000):
97         r = m2[i].tolist()
98         r.append(self.test_data[i][-1])
99         test_proj.append(r)
100    self.test_data = test_proj
101    def getDistance(self, data1, data2):
102        dist = 0.0
103        # calculate the distance between two data
104        d1 = numpy.array(data1[:-1])
105        d2 = numpy.array(data2[:-1])
106        dist = numpy.linalg.norm(d1 - d2)
107        return dist
108    def getKNeighbors(self, k, test_example):
109        distances = []
110        neighbors = []
111        # calculate the distances from the test example to all the training data
112        for i in self.training_data:
113            distance = self.getDistance(test_example, i)
114            distances.append((distance, i))
115        # sort all the distances
116        distances.sort()
117        # get the k closest neighbors
118        for i in range(0, k):
119            neighbors.append(distances[i][1])
120        return neighbors
```

```
121 def getPrediction(self, neighbors):
122     # dictionary that maps the label to the number of times this label appears
123     labels = {}
124     # initialize the dictionary
125     for i in neighbors:
126         labels[i[-1]] = 0
127     # update the dictionary
128     for i in neighbors:
129         labels[i[-1]] = labels[i[-1]] + 1
130     max_count = 0
131     # predict the majority
132     for i in labels:
133         if labels[i] > max_count:
134             max_count = labels[i]
135     # break tie randomly
136     predictions = []
137     for i in labels:
138         if labels[i] == max_count:
139             predictions.append(i)
140     index = random.randint(0, len(predictions) - 1)
141     return predictions[index]
142 def getError(self, predictions, test_examples):
143     if len(predictions) != len(test_examples):
144         print("Size does not match")
145         return
146     # get the errors
147     errors = 0.0
148     for i in range(0, len(predictions)):
149         if predictions[i] != test_examples[i][-1]:
150             errors = errors + 1.0
151     errors = errors / (float(len(test_examples)))
152     return errors
```

```
154 if __name__ == '__main__':
155     # create the solver
156     solver = Solver()
157     # load the data
158     print("loading data")
159     solver.load()
160     print("projecting data")
161     # project the data onto the projection matrix
162     solver.project()
163     # list that contains all the predictions
164     predictions = []
165     # calculate the training error
166     for i in solver.training_data:
167         # get the k nearest neighbors
168         #kNeighbors = solver.getKNeighbors(1, i)
169         #kNeighbors = solver.getKNeighbors(5, i)
170         #kNeighbors = solver.getKNeighbors(9, i)
171         kNeighbors = solver.getKNeighbors(15, i)
172         # get the prediction
173         prediction = solver.getPrediction(kNeighbors)
174         # add the prediction to the prediction list
175         predictions.append(prediction)
176     # get the number of errors
177     training_error = solver.getError(predictions, solver.training_data)
178     # print out the errors
179     print("Training error ", training_error)
180     # calculate the validation error
181     for i in solver.validation_data:
182         # get the k nearest neighbors
183         #kNeighbors = solver.getKNeighbors(1, i)
184         #kNeighbors = solver.getKNeighbors(5, i)
185         #kNeighbors = solver.getKNeighbors(9, i)
186         kNeighbors = solver.getKNeighbors(15, i)
187         # get the prediction
188         prediction = solver.getPrediction(kNeighbors)
189         # add the prediction to the prediction list
190         predictions.append(prediction)
191     # get the number of errors
192     validation_error = solver.getError(predictions, solver.validation_data)
```

```
193     # print out the errors
194     print("Validation error ", validation_error)
195     # calculate the test error
196     for i in solver.test_data:
197         # get the k nearest neighbors
198         kNeighbors = solver.getKNeighbors(15, i)
199         # get the prediction
200         prediction = solver.getPrediction(kNeighbors)
201         # add the prediction to the prediciton list
202         predictions.append(prediction)
203     # get the number of errors
204     test_error = solver.getError(predictions, solver.test_data)
205     # print out the errors
206     print("Validation error ", test_error)
```