
Image Classification with the K-Nearest-Neighbors Algorithm

Jacky Dam
University of California San Diego
San Diego, California
jdam@ucsd.edu

Abstract

In this assignment, I utilized python to create an algorithm based on k-nearest-neighbors to classify data sets containing images of numbers and the correct label for each number. What makes this algorithm special is the fact that it classifies images of numbers from a test set that the algorithm has never seen before by learning distinct characteristics from the training set. After implementing this algorithm in python, I found that the algorithm performed adequately with an error of about 9% evaluated on the test set. Additionally, I found that after performing dimensionality reduction using projections which drastically reduced runtime, the algorithm performed with an error of about 29% evaluated on the test set.

1 Data Set:

<http://cseweb.ucsd.edu/classes/wi21/cse151A-a/homeworks.html>

2 Evaluation:

Classifier:	k	Training Error:	Validation Error:
C_1	1	0	0.082
C_5	5	0.055	0.099
C_9	9	0.068	0.105
C_{15}	15	0.093	0.106

Test Error of C_1 : 0.094

3 Evaluation w/ Projection:

Classifier:	k	Training Error:	Validation Error:
C_1	1	0	0.32
C_5	5	0.1905	0.296
C_9	9	0.2285	0.29
C_{15}	15	0.2565	0.291

Test Error of C_9 : 0.29

4 Summary:

After implementing the algorithm in part 1, I found that C_1 performed the best on the validation set with an error of about 8%. I then utilized this classifier to evaluate the algorithm using the test set and found that with the test set, the algorithm performed with an error of about 9%.

In part 2, I utilized the same algorithm but utilized dimensionality reduction by projecting each image onto some projection matrix. This drastically reduced the runtime of the code by more than half, but I found that the accuracy of the algorithm was greatly reduced as the algorithm had an error of about 29% on the test set. Although the accuracy of the algorithm was drastically reduced by dimensionality reduction, it is evident that there is a trade off between computing power and performance.

```

1 import numpy as np
2
3 def extract_data(file):
4     data = []
5     with open(file) as f:
6         for line in f.readlines():
7             data.append([int(i) for i in line.strip().split(" ")])
8
9     return data
10
11 def extract_double_data(file):
12     data = []
13     with open(file) as f:
14         for line in f.readlines():
15             data.append([float(i) for i in line.strip().split(" ")])
16
17     return data
18
19 def find_distance(vector1, vector2):
20     vector = np.square(np.subtract(vector1, vector2))
21     distance = np.sqrt(np.sum(vector))
22
23     return distance
24
25 def find_neighbors(train_set, input_vector, k):
26     distance_set = []
27     neighbors_set = []
28     for train_vector in train_set:
29         distance = find_distance(train_vector[:784], input_vector)
30         distance_set.append((train_vector, distance))
31     distance_set.sort(key=lambda tup: tup[1])
32
33     for i in range(k):
34         neighbors_set.append(distance_set[i][0])
35
36     return neighbors_set
37
38 train_set = extract_data("pa1train.txt")
39 val_set = extract_data("pa1validate.txt")
40 test_set = extract_data("pa1test.txt")
41 proj_matrix = extract_double_data("projection.txt")
42
43 # This chunk of code evaluates the error of the KNN algorithm (part 1).
44
45 output_labels = []
46 actual_labels = []
47 error_count = 0
48
49 for j in range(len(test_set)):
50     neighbors = find_neighbors(train_set, test_set[j][:784], 1)
51     output = [int(i[-1]) for i in neighbors]
52     label = max(set(output), key=output.count)
53
54     if label != test_set[j][-1]:
55         error_count += 1
56
57 print(error_count/len(test_set))

```

```
58
59
60 # This chunk of code utilizes projections to evaluate the error of the KNN algorithm (part 2).
61
62 train_label = [i[-1] for i in train_set]
63 val_label = [i[-1] for i in val_set]
64 test_label = [i[-1] for i in test_set]
65
66 train_set = np.column_stack((np.matmul(np.dot([i[:784] for i in train_set], proj_matrix),
67 np.transpose(proj_matrix))/np.linalg.norm(proj_matrix), train_label))
68 val_set = np.column_stack((np.matmul(np.dot([i[:784] for i in val_set], proj_matrix),
69 np.transpose(proj_matrix))/np.linalg.norm(proj_matrix), val_label))
70 test_set = np.column_stack((np.matmul(np.dot([i[:784] for i in test_set], proj_matrix),
71 np.transpose(proj_matrix))/np.linalg.norm(proj_matrix), test_label))
72
73
74 output_labels = []
75 actual_labels = []
76 error_count = 0
77
78 for j in range(len(test_set)):
79     neighbors = find_neighbors(train_set, test_set[j][:784], 9)
80     output = [int(i[-1]) for i in neighbors]
81     label = int(max(set(output), key=output.count))
82
83     if label != test_set[j][-1]:
84         error_count += 1
85
86 print(error_count/len(test_set))
87
```