

k-nearest-neighbors

Problem

The goal of this program is to build a model that implements the k-nearest-neighbors algorithm. This model will either be a classifier or regression model, depending on user specification, that will receive training data on which to train the model, and testing data which will predict the class or likely class for each test sample provided to the predictor.

Development Environment

This program was developed on two platforms with the Python 3.7 library, though the program was mostly tested using the CuPy instead of NumPy package. This development environment consisted of the PyCharm Ultimate IDE on a CUDA-enabled GPU (GeForce GTX 1050). The other environment did not have a dedicated graphics card, which seriously hindered the testing time and allowed only a limited number of test cases, which are not presented in this paper.

Algorithms implemented

The algorithm used for the model's predictor function is the k-nearest-neighbors. The implementation was designed for the following 4 user choices:

1. Weighted classifier
2. Un-weighted classifier
3. Weighted regression
4. Un-weighted regression

Therefore, the predictor first checks whether a regression or classification model will be used, and then checks if the selected model will calculate results using weights or not. The implemented algorithm for the KNN classifier model is as follows:

```
classifier(test_data): guess
```

1. Create a list for the guesses for each of the test samples
2. Create a list for the labels for each of the possible classifications
3. For every test sample x
 - a. Create a list for the votes of the label for each neighbor
 - b. For every possible classification l
 - i. Find the delta sum of the neighbors whose labels are equal to l
 - ii. If we are using weights, multiple this value by the sum of the weights of the neighbors
 - c. Find the largest of these votes to be the guess for this test sample

The implemented algorithm for the KNN regression model is as follows:

```
regression(test_data): guess
```

1. Create a list of the means of neighbors for each test sample
2. For every test sample x
 - a. Save the sum of the weights multiplied by labels of the k nearest neighbors
 - b. If we are using weights, divide this sum by the sum of the k nearest neighbors' weights
 - c. If we are not using weights, divide this sum by k

Experimental results

As expected, the running times for regression . Two datasets were used to train and test the model: the MNIST dataset and the Solar Particle dataset. Training results' times are not reported as the `fit()` function only stores the training data for the `predict()` function to later use. The results for MNIST dataset are presented in Figures 1 and 2. The tests presented from these results use the classification and regression algorithms, with both weighted and unweighted models. All results presented included 5000 training samples and up to 500 test cases.

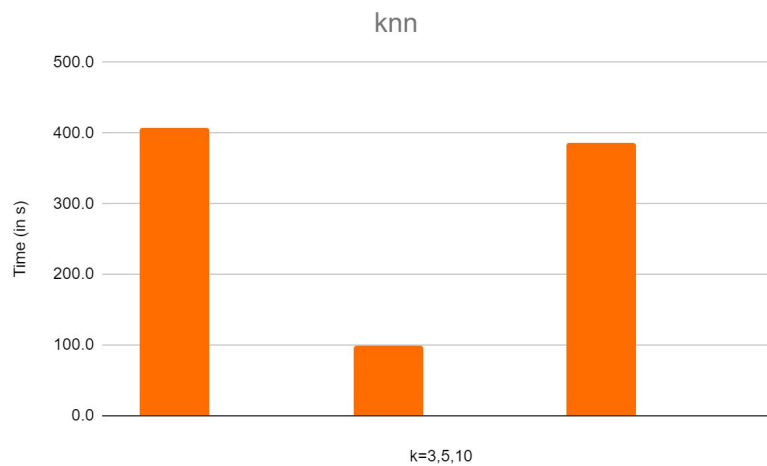


Fig 1: Average times for MNIST tests when $k = 3, 5, 10$

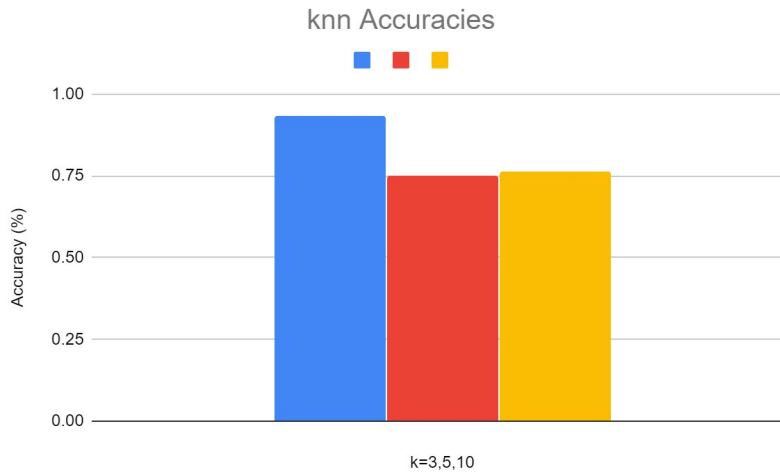


Fig 2: Average accuracies for MNIST tests when k = 3, 5, 10

Below, in Figures 3 and 4, I present the results for the Solar Particle dataset using k values of 3 and 5. Again, the tests presented from these results use the classification and regression algorithms, with both weighted and unweighted models. All results presented included 5000 training samples and up to 500 test cases.

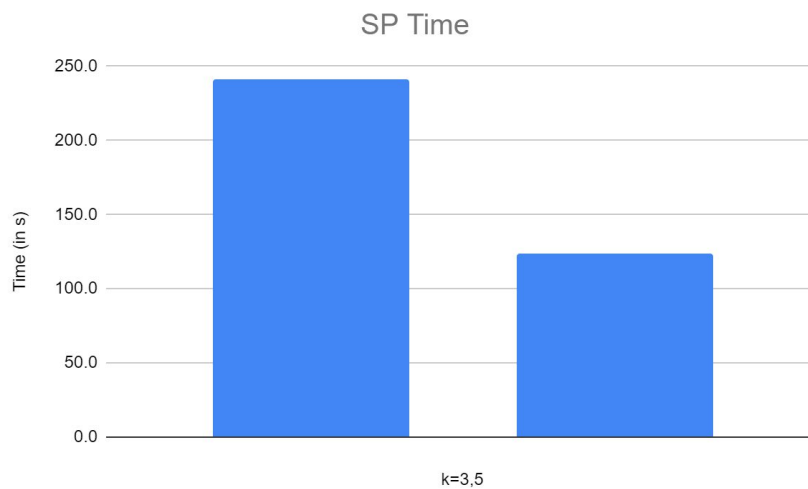


Fig 3: Average times for Solar Particle tests when k = 3, 5

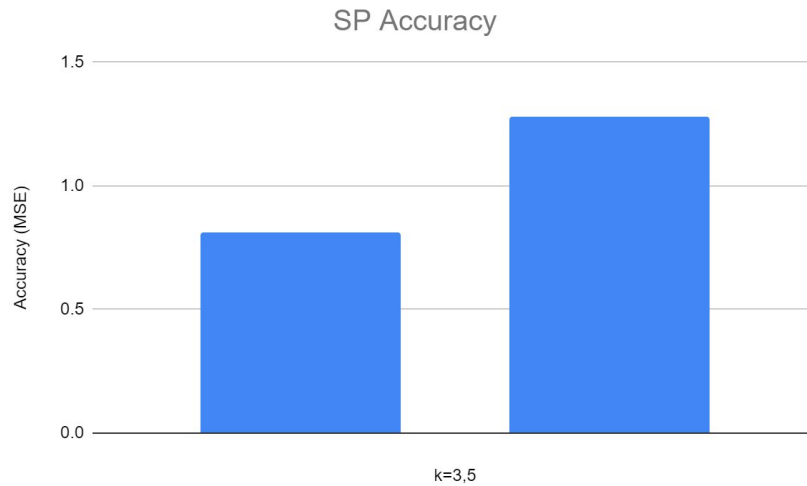


Fig 4: Average accuracies for Solar Particle tests when k = 3, 5

Optimization

In order to improve the efficiency of the program, I implemented an algorithm for attribute selection. The algorithm is as follows:

`attr_sel(train): void`

1. Find the means of each of the features across all training samples
2. Remove those features whose means are less than 0.25 of the average of features

On average, the number of features decreased by 390-400 for MNIST and decreased by 40 in the Solar Particles dataset. The time decreased but not significantly. I expected a decrease in time if the number of features decreased significantly, as reflected in the MNIST dataset, and thankfully the optimization does improve when preprocessing the data.

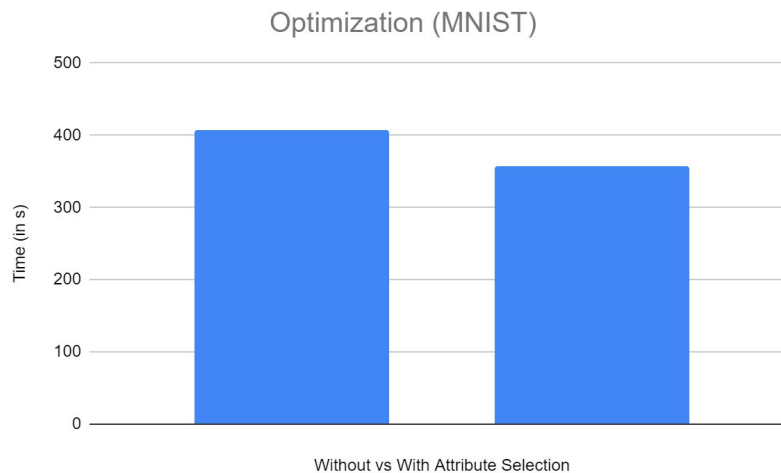


Figure 5: MNIST results with attribute selection implemented

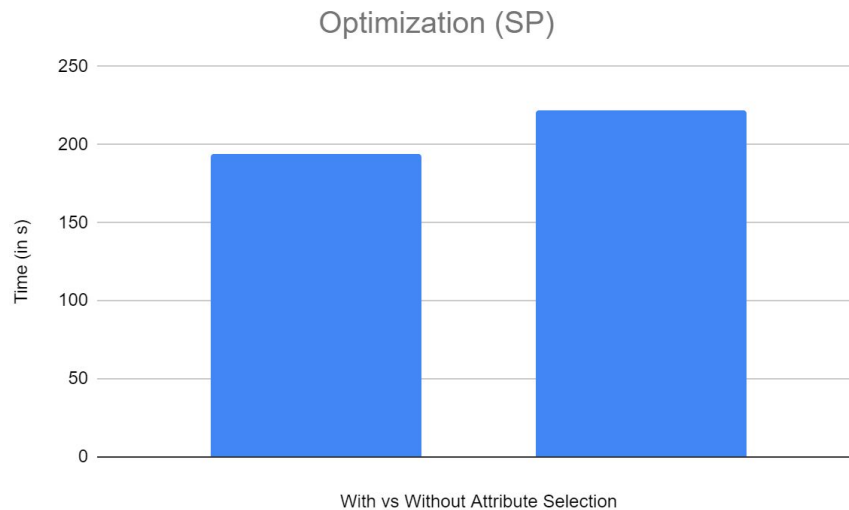


Figure 6: SP results with attribute selection implemented

For the MNIST dataset, on average, the number of features removed was 395 (out of 784). On average, the number of features removed from the Solar Particle dataset was 40 (out of 50). Since the use of attribute selection did not affect accuracy significantly, give or take 3%, it would be beneficial to implement a better algorithm as well as other optimization approaches in addition to the one used in this program.

Discussion of results

There are some limitations to my program. Firstly, my implementation may not be taking full advantage of the features of Python array slicing, and this may explain why my results do not reflect a significant decrease in time despite my implementation of optimization.

Another concern lies in my optimization algorithm, which used a fixed value for the threshold for feature variance. For the purposes of this assignment, I kept the features whose means were less than 0.25 of the average of all of the features. This, however, is dependent on the training data's features, meaning, which features are dropped will vary greatly. For smaller test cases, this did not affect accuracy but did improve efficiency. For larger test cases and varying values of k , accuracy did decrease and the time efficiency improved.

Appendix

The appendix includes the source code of the described program implemented in Python 3.7. Methods/Files authored by me are denoted with an author tag (@author mmafr). The code is also available at https://github.com/mahdafr/19w_cs5361-labs on GitHub.

```

1  import mnist
2  import cupy as np
3  import time
4  import math
5
6
7  class knn(object):
8      # Constructor
9      def __init__(self, k=3, weighted=True, classify=True):
10         self.k = k
11         self.weighted = weighted
12         self.classify = classify
13         # @author mmafr
14         self.n_matrix = None    # the list of the distances
15         self.n_weights = None
16         self.n_labels = None
17
18     # train the model
19     def fit(self, x, y):
20         self.x_train = x.astype(np.float32)
21         self.y_train = y
22
23     # test the model
24     # @author mmafr
25     def predict(self, test):
26         s = ' (weighted)' if self.weighted else ''
27         print('Classification' + s if self.classify else 'Regression' + s)
28         # test = self.attr_selection(test) # for optimization
29         self.buildNeighborList(test)
30         if self.classify:
31             return self.classifier(test)
32         else:
33             return self.regression(test)
34
35     # KNN classifier: weighted or unweighted?
36     # @author mmafr
37     def classifier(self, test):
38         guess = np.zeros(test.shape[0])
39         label = np.unique(self.y_train)    # every possible class
40         for x in range(test.shape[0]):
41             votes = np.zeros(label.shape[0]) # sums up to k
42             for l in range(label.shape[0]): # for each label
43                 d = self.delta(l, self.n_labels[x])
44                 votes[l] = d*np.sum(self.n_weights[x])
45                 votes = np.around(votes, 3)
46             guess[x] = np.argmax(votes)
47             # print("Prediction for Test=" + str(x) + " is " + str(guess[x]))
48         return guess
49
50     # @author mmafr
51     def delta(self, a, b):
52         return (b==a).sum()
53
54     # KNN regression: weighted or unweighted?
55     # @author mmafr
56     def regression(self, test):
57         mean = np.zeros(test.shape[0])
58         for x in range(test.shape[0]): # for each test example
59             mean[x] = np.sum(self.n_weights[x]*self.n_labels[x])
60             if self.weighted: # divide by sum of the weights
61                 mean[x] = mean[x]/np.sum(self.n_weights[x])
62             else: # divide by k
63                 mean[x] = mean[x]/self.k
64             # print("Prediction for Test=" + str(x) + " is " + str(mean[x]))
65         mean = np.around(mean, 0) # round for picking a class
66         return mean

```

```

67
68 # calculate the distances from each test point to other points
69 # @author mmafr
70 def buildNeighborList(self, test):
71     n_tmp = np.zeros(shape=(test.shape[0], self.x_train.shape[0]))
72     for x1 in range(test.shape[0]): # foreach test sample
73         for x2 in range(self.x_train.shape[0]): # for each train point
74             n_tmp[x1][x2] = (self.distance(test[x1], self.x_train[x2]))
75     srted = np.argsort(n_tmp, axis=1) # sort list of neighbors
76     # build the knn matrix and weights
77     self.k_neighbors(n_tmp, srted)
78     self.weights()
79     print("Built neighbor matrix.")
80
81 # calculate the euclidean distance between point x1 and x2
82 # @author mmafr
83 def distance(self, x1, x2):
84     d = 0
85     for f in range(len(x1)):
86         d += pow((x1[f] - x2[f]), 2)
87     return math.sqrt(d)
88
89 # determines the list of the k nearest-neighbors
90 # @author mmafr
91 def k_neighbors(self, neighbors, nearest):
92     # each test sample has k nearest neighbors
93     self.n_matrix = np.zeros(shape=(neighbors.shape[0], self.k))
94     self.n_labels = np.zeros(self.n_matrix.shape)
95     # self.attr_selection() # preprocessing: optimize time
96     for i in range(self.n_matrix.shape[0]):
97         for j in range(self.k):
98             self.n_matrix[i][j] = neighbors[i][nearest[i][j]]
99             self.n_labels[i][j] = self.y_train[int(nearest[i][j])]
100
101 # drop the features with the lowest variance
102 # @author mmafr
103 def attr_selection(self, test):
104     start = self.x_train.shape[1]
105     mean = np.mean(self.x_train, axis=0) # means of features
106     mask = (self.x_train > 0.25*np.mean(mean))
107     self.x_train = np.asarray(self.x_train[:, mask.any(axis=0)])
108     x_test = test[:, mask.any(axis=0)]
109     print("Dropped " + str(start - self.x_train.shape[1]) + " features.")
110     return x_test
111
112 # get the weights, list of ones if not weighted
113 # @author mmafr
114 def weights(self):
115     if self.classify:
116         self.n_weights = np.zeros(self.n_matrix.shape)
117         self.n_weights[:, 0] = 1 # all others are 0
118     else:
119         self.n_weights = np.ones(self.n_matrix.shape)
120     if self.weighted: # if weighted, w=1/d(x,xi)
121         self.n_weights = 1/self.n_matrix
122     # mask = np.isin(self.sorted, self.neighbor)
123     # weight = 1/weight[mask]
124
125
126 if __name__ == "__main__":
127     TESTS = 10
128     TRAIN = 1000
129     THEK = 3
130     print('MNIST dataset')
131     x_train, y_train, x_test, y_test = mnist.load()
132     x_train = x_train[:TRAIN:]

```

```

133 y_train = y_train[:TRAIN:]
134 x_test = x_test[:TESTS:]
135 y_test = y_test[:TESTS:]
136 print('Training size = ' + str(x_train.shape[0]))
137 print('Testing size = ' + str(x_test.shape[0]))
138 print("K = " + str(THEK))
139 model = knn(classify=True, weighted=False)
140
141 start = time.time()
142 model.fit(x_train, y_train)
143 elapsed_time = time.time()-start
144 print('Elapsed_time training {0:.6f} '.format(elapsed_time))
145
146 start = time.time()
147 pred = model.predict(x_test)
148 elapsed_time = time.time()-start
149 print('Elapsed_time testing {0:.6f} '.format(elapsed_time))
150
151 y_test = np.asarray(y_test) # for CuPy
152 print('Accuracy:', np.sum(pred == y_test)/len(y_test))
153
154 # print('\nSolar particle dataset')
155 # dir = 'D:\Google Drive\skool\CS 5361\datasets\lab1\'
156 # x_train = np.load(dir + 'x_ray_data_train.npy')
157 # y_train = np.load(dir + 'x_ray_target_train.npy')
158 # x_test = np.load(dir + 'x_ray_data_test.npy')
159 # y_test = np.load(dir + 'x_ray_target_test.npy')
160 # y_test = np.asarray(y_test) # for CuPy
161 # x_train = x_train[:TRAIN]
162 # y_train = y_train[:TRAIN]
163 # x_test = x_test[:TESTS:]
164 # y_test = y_test[:TESTS:]
165 # print('Training size = ' + str(x_train.shape[0]))
166 # print('Testing size = ' + str(x_test.shape[0]))
167 # print("K = " + str(THEK))
168 # model = knn(k=THEK, classify=True, weighted=True)
169 #
170 # start = time.time()
171 # model.fit(x_train, y_train)
172 # elapsed_time = time.time()-start
173 # print('Elapsed_time training {0:.6f} '.format(elapsed_time))
174 #
175 # start = time.time()
176 # pred = model.predict(x_test)
177 # elapsed_time = time.time()-start
178 # print('Elapsed_time testing {0:.6f} '.format(elapsed_time))
179 #
180 # print('Mean square error:', np.mean(np.square(pred-y_test)))
181

```