

ROB313 Intro to Learning from Data - Assignment 1

Objectives

The objective of this assignment is to familiarize one with the applications and varying implementations of the k-NN algorithm. In this respect, the k-NN algorithm will be used for both regression and classification of several datasets. Additionally, various parameters of the algorithm will be simulated to analyze the effect on the k-NN's performance (e.g. distance metrics, feature space dimensions). In doing so, one will also become familiar with the available linear algebra tools/libraries for various computations. The performance of the k-NN algorithm will be compared to the Linear Regression Model for deeper understanding. Altogether, this assignment serves to provide a first-hand experience in applying theoretical machine learning knowledge by implementing two introductory level algorithms

Code Structure and Strategies

The code structure was designed to be modular and easily understandable. To run the code for each question, the only modifications needed to be made to the file is to uncomment the corresponding section outlined in the main block. The remaining code is structured as follows:

*Q1: **ff_regression**(data) – 5-fold regression function, **ff_regression_test**(data) – regression test*

*Q2: **one_fold_classification**(data), **classification_test**(data) – classification validation and test*

Questions 1 and 2 each share a common function called **run_model**(model_type, dataset) from which the Q1 and Q2 functions are called from, depending on the *model_type* parameter. The *run_model* function is called directly from the main block.

*Q3: **regression_performance**(data, method) – Modified k-NN algorithms*

This function is called many times to compute the Test RMSE for the provided dataset. However, the *method* parameter specifies what type of computation should take place (e.g. *method* = 'a' will compute predictions using the double loop over test points). This function is called directly from the main block.

*Q4: **linear_regression**(data, model_type) – Linear regression function*

This function is called directly from the main block and will compute predictions and Test RMSE/Test Accuracy Ratio using a linear model for either regression or classification datasets. The parameter, *model_type*, specified whether to run classification or regression.

The first strategy deployed was to make the code very modular. Thus, easily interchangeable functions are run directly from the main block for each question. The functions were designed to be generalizable so that they can be re-used for datasets with differing feature space dimensions. Furthermore, the final code has gone through several iterations to optimize the runtime of each function. Originally, performing regression and classification for the larger datasets would take several hours at a time. Changes such as inverting the loop structure of the function significantly decreased the overall runtime. For example, in Q1 and Q2, the distances between a given test

point and all training points are computed only once, and then used to obtain the k-nearest neighbours for many of k values. Finally, print statements are neatly organized throughout the main block so that the results can be easily obtained by the operator when executing the code. Neatness and organization aided in the debugging process so that little time was wasted searching for incorrect results.

Q1 - k-NN Algorithm for Regression

1-a) Implement k-NN Algorithm for Regression Datasets and use 5-fold Cross Validation:

The results of running the k-NN regression on all regression datasets are provided below. A *Cross-Validation RMSE*, *Test RMSE*, *Estimated k*, and *preferred Distance Metric* is reported for each dataset. Note that for each dataset, k values on a range (1, 30) were simulated. From these results, we can see a correspondence between the dimensionality of the feature space to the estimated k value. Additionally, we see that datasets with close feature space dimensions prefer the same distance metric.

Table 1 – k-NN Regression Results

Dataset	Estimated k	Distance Metric	CV RMSE	Test RMSE
mauna_loa	2	L2-norm	0.0318041026269	0.440704890355
rosenbrock	2	L2-norm	0.330741584025	0.247598644133
pumadyn32nm	25	L-inf-norm	0.87221127042	0.832480142444

1-b) Cross-Validation Loss Curve and Prediction on Test Set for Mauna Loa (L2-norm)

Figure 1 illustrates the cross-validation loss (merging predictions from all splits) across increasing values of k for the Mauna Loa dataset. As expected, we see that the loss is minimized for a k value of 2. This corresponds to the Mauna Loa result shown in *Table 1*. Further, we see that the loss rapidly increases for higher k values which suggests that computing an estimate with more neighbours negatively effects the accuracy of predictions.

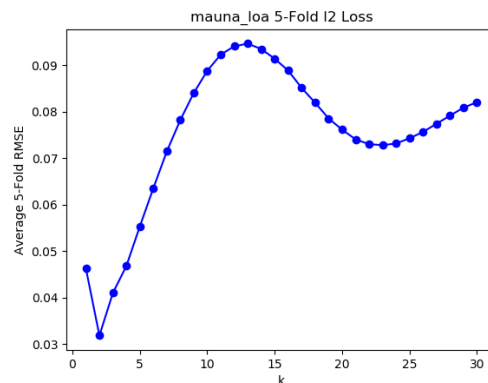


Figure 1 – Cross Validation Loss Curve for Mauna Loa

In *Figure 2*, we have the prediction of the k-NN regression on the test set for Mauna Loa. Notice that for all test points, the predicted values are identical. Intuitively, this makes sense as the Mauna Loa dataset represents data taken over *time*. The 1-dimensional feature space of time values for the test set contains data points that were collected after all the data points in the training and validation set, chronologically. Thus, the k nearest neighbours for all test points are the same, yielding identical predictions. It is for this reason that the Test RMSE shown in *Table 1* is significantly larger than the Cross-Validation RMSE for this particular dataset.

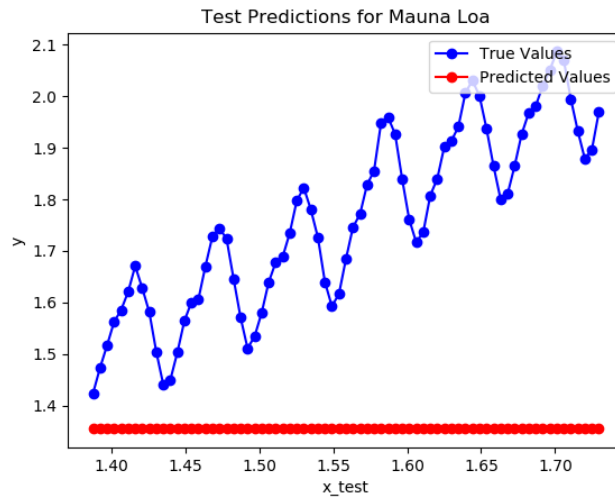


Figure 2 – Prediction on Test Set for Mauna Loa

Q2 - k-NN Algorithm for Classification

A similar k-NN algorithm was used as a classification model and was applied to two classification datasets. The results are summarized in *Table 2*, containing the *Estimated k* , *preferred Distance Metric*, *Validation Accuracy Ratio*, and the *Test Accuracy Ratio* for each dataset.

Table 2 – k-NN Classification Results

Dataset	Estimated k	Distance Metric	Validation Ratio	Test Ratio
iris	15	L1-norm	0.9032258064516129	1.0
mnist_small	1	L2-norm	0.95	0.959

Q3 - k-NN Performance with Varying Computational Modifications

Four versions of the k-NN algorithm were implemented to compute predictions on the rosenbrock ($k=5$, $n_{\text{train}} = 5000$) regression test set over increasing values of d (dimension of feature space). The four versions were: *brute force double loop* (a), *half-vectorized method* (b), *fully-vectorized method* (c), *k-d tree method* (d). Using the L2-norm to compute distances, the

performance metrics *Run-time and Test RMSE*, were tracked over d values on a range (2, 9). The quantities are summarized in *Table 3* and displayed in *Figures 3-4*.

Table 3 – k-NN Performance Results with Various Computational Modifications (plotted on Figures 3-4)

d	Methods (a, b, c, d)							
	Double Loop (a)		Half-Vectorized (b)		Fully-Vectorized (c)		k-d Tree (d)	
---	Time [s]	RMSE	Time [s]	RMSE	Time [s]	RMSE	Time [s]	RMSE
2	172.363	0.267	2.505	0.267	0.157	0.267	0.003	0.267
3	106.918	0.379	2.873	0.379	0.155	0.379	0.005	0.379
4	109.243	0.420	2.899	0.420	0.171	0.420	0.007	0.420
5	305.698	0.521	3.007	0.521	0.158	0.521	0.011	0.521
6	290.247	0.611	3.012	0.611	0.164	0.611	0.018	0.611
7	129.156	0.689	3.485	0.689	0.170	0.689	0.029	0.689
8	310.322	0.747	6.338	0.747	0.204	0.747	0.174	0.747
9	153.278	0.800	3.904	0.800	0.180	0.800	0.064	0.800

Figure 3 depicts the prediction runtimes of all four methods over increasing feature space dimension. As we can see, the Double Loop method take significantly longer than the vectorized and k-d tree methods. Analyzing *Table 3*, we can see that the k-d tree implementation is consistently the best performer in terms of time, followed by the fully-vectorized and half-vectorized methods, although the k-d tree and fully-vectorized methods come quite close. An interesting observation is that the double loop method seems to be extremely sensitive to the size of the feature space, while the other three methods stay quite consistent in runtimes. Overall, vectorizing the code significantly reduced the runtime of the algorithm, and using a k-d data structure further improves the computational performance in comparison to the brute-force method.

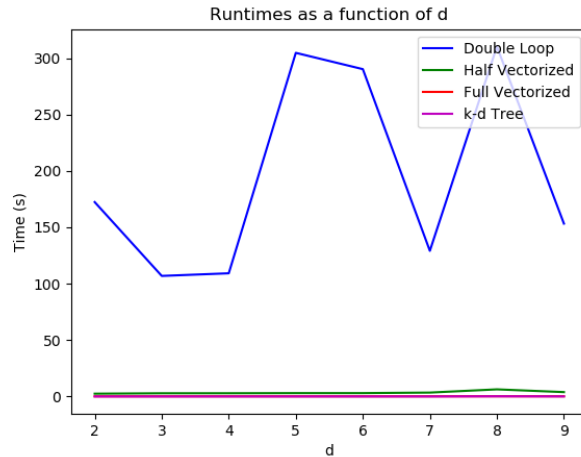


Figure 3 – Prediction Runtimes for all k-NN Methods

Figure 4 displays the Test RMSE values for all methods over the increasing feature space dimension. As expected, all versions of the k-NN algorithm predict the same values, and hence their test RMSE values should be identical. This corresponds to the overlapping curves below.

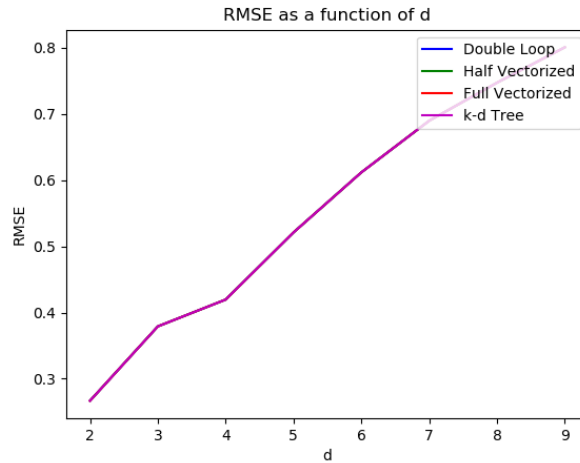


Figure 4 – Test RMSE for all k-NN Methods

Q4 - Linear Regression Model for Regression and Classification

Two versions of a linear regression model were implemented for regression and classification of all datasets. The test RMSE values (regression) and test accuracy ratios (classification) are reported in *Table 4*. Comparing the results of the linear model to the k-NN algorithm we see several things. For the regression sets, the test RMSE values suffered a significant increase for rosenbrock and a minor increase pumadyn32nm. However, the test RMSE value decreased for the Mauna Loa dataset. The result for Mauna Loa is expected, as the linear regression model is able to better fit the data with an increasing slope, which the k-NN algorithm was unable to do (see *Figure 2*). For classification, the linear regression model slight underperformed in comparison to the k-NN algorithm for both iris and mnist_small datasets. Thus, the k-NN algorithm outperforms the linear regression model over all datasets except for Mauna Loa.

Table 4 – Linear Regression Model Test Results

Dataset	Test RMSE	Test Accuracy Ratio
mauna_loa	0.349388310499	---
rosenbrock	0.984087203069	---
pumadyn32nm	0.86225124366	---
iris	---	0.866666666667
mnist_small	---	0.855

```

1 import numpy as np
2 import time
3 import math
4 from matplotlib import pyplot as plt
5 from data_utils import load_dataset
6 import heapq
7 from sklearn import neighbors
8
9 __author__ = 'Christopher Agia (1003243509)'
10 __date__ = 'February 12, 2019'
11
12
13 # Root Mean Squared Error Function
14 def rmse(y_test, y_estimates):
15     return np.sqrt(np.average((y_test-y_estimates)**2))
16
17 # Norm Utility Functions
18 def l1_norm(x1, x2):
19     return np.linalg.norm([x1-x2], ord=1)
20
21
22 def l2_norm(x1, x2):
23     return np.linalg.norm([x1-x2], ord=2)
24
25
26 def linf_norm(x1, x2):
27     return np.linalg.norm([x1-x2], ord=np.inf)
28
29
30 def ff_regression(x_train, x_valid, y_train, y_valid, distance_functions, k_list=None):
31
32     assert((len(x_train)+len(x_valid)) == (len(y_train)+len(y_valid)))
33
34     rmse_vals = {}
35     error = []
36
37     x_total = np.vstack([x_train, x_valid])
38     y_total = np.vstack([y_train, y_valid])
39     np.random.seed(5)
40     np.random.shuffle(x_total)
41     np.random.seed(5)
42     np.random.shuffle(y_total)
43
44     ff_length = len(x_total)//5
45
46     # Trial K values to try
47     if not k_list:
48         k_list = list(range(0, 30))
49
50     # Iterate over 5-folds
51     for i in range(5):
52
53         # Obtain validation and train set associated with fold i
54         y_valid = y_total[i * ff_length:(i + 1) * ff_length]
55         y_train = np.vstack([y_total[:i * ff_length], y_total[(i + 1) * ff_length:]])
56         x_valid = x_total[i * ff_length:(i + 1) * ff_length]
57         x_train = np.vstack([x_total[:i * ff_length], x_total[(i + 1) * ff_length:]])
58
59         # Loop over distance functions (e.g. l1_norm, l2_norm, linf_norm)
60         for func in distance_functions:
61             y_est = {}
62             # Compute distances according to distance function for one validation point
63             for j in range(ff_length):
64                 d = []
65                 for t in range(len(x_train)):
66                     d.append((func(x_train[t], x_valid[j]), y_train[t]))
67
68                 # Sort distances to take the nearest k-values
69                 d.sort(key=lambda x: x[0])
70
71                 # Compute y_estimate for validation point j
72                 for k in k_list:
73                     y = 0
74                     for elem in d[:k+1]:
75                         y += elem[1]
76
77                     if k not in y_est:
78                         y_est[k] = []
79

```

```

80     y_est[k].append(y/(k+1))
81
82     # Compute root mean squared error for each k-value
83     for k in k_list:
84         if (func, k) not in rmse_vals:
85             rmse_vals[(func, k)] = []
86             rmse_vals[(func, k)].append(rmse(y_valid, y_est[k]))
87
88
89     # Error will contain 30 rmse-error values as per k for each function
90     for func, k in rmse_vals:
91         ff_error = sum(rmse_vals[(func, k)]) / 5
92         error.append((k+1, func, ff_error))
93
94     return error
95
96
97 def ff_regression_test(x_train, x_valid, x_test, y_train, y_valid, y_test, k, func, plot=False):
98
99     # Create total training set (no validation)
100    x_total = np.vstack([x_train, x_valid])
101    y_total = np.vstack([y_train, y_valid])
102
103    # Predictions for each test data point will be stored in this list
104    predictions = []
105
106    # Iterate over test points
107    for elem in x_test:
108
109        # Compute distances between all training points and test point
110        d = []
111        for i in range(len(x_total)):
112            d.append((func(elem, x_total[i]), y_total[i]))
113
114        # Sort in terms of increasing distance
115        d.sort(key=lambda x: x[0])
116
117        # Average over k nearest y values
118        y_est = 0
119        for item in d[:k]:
120            y_est += item[1]
121        avg = y_est/k
122
123        # Append to predictions
124        predictions.append(avg)
125
126    test_error = rmse(y_test, predictions)
127
128    if plot:
129        # Predictions on test set
130        plt.figure(2)
131        plt.plot(x_test, y_test, '-bo', label='True Values')
132        plt.plot(x_test, predictions, '-ro', label='Predicted Values')
133        plt.title('Test Predictions for Mauna Loa')
134        plt.xlabel('x_test')
135        plt.ylabel('y')
136        plt.legend(loc='upper right')
137        plt.savefig('mauna_loa_prediction.png')
138
139    return test_error
140
141
142 def one_fold_classification(x_train, y_train, x_valid, y_valid, distance_functions, k_list=None):
143
144     assert ((len(x_train) + len(x_valid)) == (len(y_train) + len(y_valid)))
145
146     tally = {}
147
148     if not k_list:
149         k_list = list(range(0, 30))
150
151     # Loop over distance functions (e.g. L1_norm, L2_norm, Linf_norm)
152     for func in distance_functions:
153
154         # Compute distances according to distance function for one validation point
155         for j in range(len(x_valid)):
156             d = []
157             for t in range(len(x_train)):
158                 d.append((func(x_train[t], x_valid[j]), y_train[t]))
159

```

```

160     # Sort distances to take the nearest k-values
161     d.sort(key=lambda x: x[0])
162
163     classes = {}
164     # Identify k-nearest neighbours for x_valid[j]
165     for k in k_list:
166         classes[k] = []
167         for elem in d[:k + 1]:
168             classes[k].append(elem[1])
169
170     for k in k_list:
171         occurance = {}
172         for point in classes[k]:
173             if str(point) not in occurance:
174                 occurance[str(point)] = (point, 0)
175                 occurance[str(point)] = (point, occurance[str(point)][1] + 1)
176
177     occur_list = list(occurance.values())
178     occur_list.sort(key=lambda x: x[1], reverse=True)
179
180     if np.all(occur_list[0][0] == y_valid[j]):
181         if (k + 1, func) not in tally:
182             tally[(k + 1, func)] = 0
183             tally[(k + 1, func)] += 1
184
185     results = []
186     for k, func in tally:
187         ratio = tally[(k, func)]/len(y_valid)
188         result = (k, func, ratio)
189         results.append(result)
190
191     return results
192
193
194 def classification_test(x_train, x_valid, x_test, y_train, y_valid, y_test, k, func):
195
196     # Create total training set (no validation)
197     x_total = np.vstack([x_train, x_valid])
198     y_total = np.vstack([y_train, y_valid])
199
200     tally = 0
201     # Iterate over test points
202     for i in range(len(x_test)):
203
204         # Compute distances between all training points and test point
205         d = []
206         for j in range(len(x_total)):
207             d.append((func(x_test[i], x_total[j]), y_total[j]))
208
209         # Sort in terms of increasing distance
210         d.sort(key=lambda x: x[0])
211
212         occurance = {}
213         for item in d[: k]:
214             if str(item[1]) not in occurance:
215                 occurance[str(item[1])] = (item[1], 0)
216                 occurance[str(item[1])] = (item[1], occurance[str(item[1])][1] + 1)
217
218         occur_list = list(occurance.values())
219         occur_list.sort(key=lambda x: x[1], reverse=True)
220
221         if np.all(occur_list[0][0] == y_test[i]):
222             tally += 1
223
224     return tally/len(x_test)
225
226
227 def run_model(model_type, dataset):
228
229     functions = [l1_norm, l2_norm, linf_norm]
230
231     if model_type == 'regression':
232
233         if dataset == 'rosenbrock':
234             x_train, x_valid, x_test, y_train, y_valid, y_test = load_dataset(dataset, n_train=1000, d=2)
235         else:
236             x_train, x_valid, x_test, y_train, y_valid, y_test = load_dataset(dataset)
237
238         if dataset == 'mauna_loa':
239             result = ff_regression(x_train, x_valid, y_train, y_valid, [l2_norm])

```



```

240
241     # Result = (k, Function, RMSE Error)
242     result.sort(key=lambda x: x[0])
243
244     k_vals = []
245     error_vals = []
246     for k, func, error in result:
247         k_vals.append(k)
248         error_vals.append(error)
249
250     plt.figure(1)
251     plt.plot(k_vals, error_vals, '-bo')
252     plt.xlabel('k')
253     plt.ylabel('Average 5-Fold RMSE')
254     plt.title(dataset + ' 5-Fold L2 Loss')
255     plt.savefig(dataset + 'l2_loss.png')
256
257     result.sort(key=lambda x: x[2])
258     k_min = result[0][0]
259     func = result[0][1]
260     test_error = ff_regression_test(x_train, x_valid, x_test, y_train, y_valid, y_test, k_min, func, plot=True)
261
262     print('-----')
263     print('Results for Mauna Loa with L2 Norm :')
264     print('')
265     print('Optimal k: ' + str(k_min))
266     print('Optimal Distance Metric: ' + str(func))
267     print('Five fold RMSE: ' + str(result[0][2]))
268     print('Test RMSE: ' + str(test_error))
269     print('')
270
271     result = ff_regression(x_train, x_valid, y_train, y_valid, functions)
272     result.sort(key=lambda x: x[2])
273     k_min = result[0][0]
274     func = result[0][1]
275     test_error = ff_regression_test(x_train, x_valid, x_test, y_train, y_valid, y_test, k_min, func)
276
277     return result[0][0], result[0][1], result[0][2], test_error
278
279 elif model_type == 'classification':
280
281     x_train, x_valid, x_test, y_train, y_valid, y_test = load_dataset(dataset)
282     one_fold_result = one_fold_classification(x_train, y_train, x_valid, y_valid, functions)
283     one_fold_result.sort(key=lambda x: x[2], reverse=True)
284     # ff_result = (k, func, ratio)
285     k_min = one_fold_result[0][0]
286     func = one_fold_result[0][1]
287     test_ratio = classification_test(x_train, x_valid, x_test, y_train, y_valid, y_test, k_min, func)
288
289     return k_min, func, one_fold_result[0][2], test_ratio
290
291 return 0
292
293
294 def regression_performance(x_total, x_test, y_total, y_test, k, f, method):
295
296     start = time.time()
297
298     # Predictions for each test data point will be stored in this List
299     predictions = []
300
301     # Brute-force Double For Loop Method
302     if method == 'a':
303
304         # Iterate over test points
305         for elem in x_test:
306
307             # Compute distances between all training points and test point
308             d = []
309             for j in range(len(x_total)):
310                 d.append((f(elem, x_total[j]), y_total[j]))
311
312             # Sort in terms of increasing distance
313             d.sort(key=lambda x: x[0])
314
315             # Average over k nearest y values
316             y_est = 0
317             for item in d[:k]:
318                 y_est += item[1]
319             avg = y_est/k

```

```

320
321     # Append to predictions
322     predictions.append(avg)
323
324     test_error = rmse(y_test, predictions)
325
326     # Half-Vectorization Method
327     elif method == 'b':
328
329         for j in range(len(x_test)):
330             d = np.sqrt(np.sum(np.square(x_total - x_test[j]), axis=1))
331             k_nb = heapq.nsmallest(k, range(len(d)), d.take)
332             predictions.append(np.array([(np.average(np.take(y_total, k_nb)))]))
333
334         test_error = rmse(y_test, predictions)
335
336     # Full Vectorization Method
337     elif method == 'c':
338
339         d = np.sqrt(-2 * np.dot(x_test, x_total.T) + np.sum(x_total ** 2, axis=1) + np.sum(x_test ** 2, axis=1)[:, np.newaxis])
340         k_nb = np.argpartition(d, kth=k, axis=1)[:k]
341         predictions = np.sum(y_total[k_nb], axis=1) / k
342
343         test_error = rmse(y_test, predictions)
344
345     # K-D Tree Method
346     elif method == 'd':
347
348         kdt = neighbors.KDTree(x_total)
349         d, k_nb = kdt.query(x_test, k=k)
350         predictions = np.sum(y_total[k_nb], axis=1) / k
351
352         test_error = rmse(y_test, predictions)
353
354     runtime = time.time() - start
355     return runtime, test_error
356
357
358 def linear_regression(x_train, x_valid, x_test, y_train, y_valid, y_test, model_type):
359
360     if model_type == 'regression':
361
362         x_total = np.vstack([x_train, x_valid])
363         y_total = np.vstack([y_train, y_valid])
364
365         # Create X matrix
366         X = np.ones((len(x_total), len(x_total[0]) + 1))
367         X[:, 1:] = x_total
368
369         # Compute SVD
370         U, S, Vh = np.linalg.svd(X)
371
372         # Invert Sigma
373         sig = np.diag(S)
374         filler = np.zeros((len(x_total)-len(S), len(S)))
375         sig_inv = np.linalg.pinv(np.vstack([sig, filler]))
376
377         # Compute weights and predictions
378         w = np.dot(Vh.T, np.dot(sig_inv, np.dot(U.T, y_total)))
379
380         X_test = np.ones((len(x_test), len(x_test[0]) + 1))
381         X_test[:, 1:] = x_test
382         predictions = np.dot(X_test, w)
383
384         result = rmse(y_test, predictions)
385
386     elif model_type == 'classification':
387
388         x_total = np.vstack([x_train, x_valid])
389         y_total = np.vstack([y_train, y_valid])
390
391         # Expand X matrix
392         X = np.ones((len(x_total), len(x_total[0]) + 1))
393         X[:, 1:] = x_total
394
395         # Convert to integer
396         #y_test = 1 * y_test
397
398         # Perform SVD
399         U, S, Vh = np.linalg.svd(X)

```

```

400
401     # Expand Sigma Matrix
402     sig = np.diag(S)
403     filler = np.zeros([len(x_total) - len(S), len(S)])
404     sig_inv = np.linalg.pinv(np.vstack([sig, filler]))
405
406     # Compute weights
407     w = np.dot(Vh.T, np.dot(sig_inv, np.dot(U.T, y_total)))
408
409     # Create Test Matrix
410     X_test = np.ones([len(x_test), len(x_test[0]) + 1])
411     X_test[:, 1:] = x_test
412
413     # find prediction accuracy
414     predictions = np.argmax(np.dot(X_test, w), axis=1)
415     y_test = np.argmax(1 * y_test, axis=1)
416
417     result = (predictions == y_test).sum() / len(y_test)
418
419     return result
420
421
422 if __name__ == '__main__':
423     # All Dataset Names
424     all_datasets = ['mauna_loa', 'rosenbrock', 'pumadyn32nm', 'iris', 'mnist_small']
425     regression_sets = ['mauna_loa', 'rosenbrock', 'pumadyn32nm']
426     classification_sets = ['iris', 'mnist_small']
427
428     # ----- Question 1 -----
429
430     # print('----- Overall Results for Question 1 -----')
431     # print('')
432     # for d_set in regression_sets:
433     #     k_min, metric_min, ff_rmse, test_rmse = run_model('regression', d_set)
434     #     print('-----')
435     #     print('Results for ' + d_set + ' :')
436     #     print('')
437     #     print('Optimal k: ' + str(k_min))
438     #     print('Optimal Distance Metric: ' + str(metric_min))
439     #     print('Five fold RMSE: ' + str(ff_rmse))
440     #     print('Test RMSE: ' + str(test_rmse))
441     #     print('')
442
443
444     # ----- Question 2 -----
445
446     # print('----- Overall Results for Question 2 -----')
447     # print('')
448     # for d_set in classification_sets:
449     #     k_min, metric_min, max_ratio, test_ratio = run_model('classification', d_set)
450     #     print('-----')
451     #     print('Results for ' + d_set + ' :')
452     #     print('')
453     #     print('Optimal k: ' + str(k_min))
454     #     print('Optimal Distance Metric: ' + str(metric_min))
455     #     print('Validation Ratio: ' + str(max_ratio))
456     #     print('Test Ratio: ' + str(test_ratio))
457     #     print('')
458
459
460     # ----- Question 3 -----
461
462     # result_table = {}
463     # result_table['Double Loop'] = []
464     # result_table['Half Vectorized'] = []
465     # result_table['Full Vectorized'] = []
466     # result_table['k-d Tree'] = []
467     #
468     # for ind_d in range(2, 10):
469     #
470     #     x_train, x_valid, x_test, y_train, y_valid, y_test = load_dataset('rosenbrock', n_train=5000, d=ind_d)
471     #     x_total = np.vstack([x_train, x_valid])
472     #     y_total = np.vstack([y_train, y_valid])
473     #
474     #     run_time, test_rmse = regression_performance(x_total, x_test, y_total, y_test, 5, l2_norm, 'a')
475     #     result_table['Double Loop'].append((ind_d, run_time, test_rmse))
476     #
477     #     run_time, test_rmse = regression_performance(x_total, x_test, y_total, y_test, 5, l2_norm, 'b')
478     #     result_table['Half Vectorized'].append((ind_d, run_time, test_rmse))
479     #

```

```

480 # run_time, test_rmse = regression_performance(x_total, x_test, y_total, y_test, 5, l2_norm, 'c')
481 # result_table['Full Vectorized'].append((ind_d, run_time, test_rmse))
482 #
483 # run_time, test_rmse = regression_performance(x_total, x_test, y_total, y_test, 5, l2_norm, 'd')
484 # result_table['k-d Tree'].append((ind_d, run_time, test_rmse))
485 #
486 # print('----- Overall Results for Question 3 -----')
487 # print('')
488 #
489 # m = list(result_table.keys())
490 # for r in range(8):
491 #     print('--- Results for d = ' + str(r+2) + ' ---')
492 #     print('')
493 #     print('Times-- ' + m[0] + ': ' + str(result_table[m[0]][r][1]) + ' ' + m[1] + ': ' + str(result_table[m[1]][r][1]))
494 #     print(' ' + m[2] + ': ' + str(result_table[m[2]][r][1]) + ' ' + m[3] + ': ' + str(result_table[m[3]][r][1]))
495 #     print('RMSEs-- ' + m[0] + ': ' + str(result_table[m[0]][r][2]) + ' ' + m[1] + ': ' + str(result_table[m[1]][r][2]))
496 #     print(' ' + m[2] + ': ' + str(result_table[m[2]][r][2]) + ' ' + m[3] + ': ' + str(result_table[m[3]][r][2]))
497 #     print('')
498 #
499 # d_arr = list(range(2, 10))
500 # plt.figure(3)
501 # plt.title('Runtimes as a function of d')
502 # plt.xlabel('d')
503 # plt.ylabel('Time (s)')
504 # plt.figure(4)
505 # plt.title('RMSE as a function of d')
506 # plt.xlabel('d')
507 # plt.ylabel('RMSE')
508 # count = 0
509 #
510 # for m in result_table:
511 #     if count == 0:
512 #         tab = '-b'
513 #     elif count == 1:
514 #         tab = '-g'
515 #     elif count == 2:
516 #         tab = '-r'
517 #     else:
518 #         tab = '-m'
519 #     count += 1
520 #
521 #     runtimes = []
522 #     rmsees = []
523 #     for i in range(len(result_table[m])):
524 #         runtimes.append(result_table[m][i][1])
525 #         rmsees.append(result_table[m][i][2])
526 #
527 #     plt.figure(3)
528 #     plt.plot(d_arr, runtimes, tab, label=m)
529 #
530 #     plt.figure(4)
531 #     plt.plot(d_arr, rmsees, tab, label=m)
532 #
533 #     plt.figure(3)
534 #     plt.legend(loc='upper right')
535 #     plt.savefig('runtimes_vs_d.png')
536 #
537 #     plt.figure(4)
538 #     plt.legend(loc='upper right')
539 #     plt.savefig('rmsees_vs_d.png')
540 #
541 # ----- Question 4 -----
542 #
543 # print('----- Overall Results for Question 4 -----')
544 # print('')
545 # for d_set in regression_sets:
546 #     if d_set == 'rosenbrock':
547 #         x_train, x_valid, x_test, y_train, y_valid, y_test = load_dataset(d_set, n_train=1000, d=2)
548 #     else:
549 #         x_train, x_valid, x_test, y_train, y_valid, y_test = load_dataset(d_set)
550 #
551 #     final_rmse = linear_regression(x_train, x_valid, x_test, y_train, y_valid, y_test, 'regression')
552 #     print('Test RMSE for ' + d_set + ': ' + str(final_rmse))
553 #
554 # for d_set in classification_sets:
555 #     x_train, x_valid, x_test, y_train, y_valid, y_test = load_dataset(d_set)

```

```
560 #  
561 # final_ratio = linear_regression(x_train, x_valid, x_test, y_train, y_valid, y_test, 'classification')  
562 # print('Test Ratio for ' + d_set + ': ' + str(final_ratio))
```