

Task 1**Ridge Regression Description:**Linear regression models –

A Linear regression model is used to predict values based on features of pre-existing data, to do this a line of best fit is calculated in the form of the equation " $Y \approx \beta_0 + \beta_1 X$ " (James et al, 2013, 61).

MSE and SSE –

Both SSE and MSE are attempts to evaluate how well a model created by a regression function fits the data and by extension how useful it is in predicting data points.

SSE stands for Sum of Squared Errors which is the sum of differences between the line of best fit created by and the data points, the data is squared so that if data points are below the regression model, causing negative error numbers, do not make the line look better by decreasing the overall sum.

MSE stands for Mean Squared Error which is the SSE divided by the number of data points used. This gives the average error.

Least squares solution –

To find the optimal result for regressions, the SSE or RMSE must be reduced as much as possible, giving the least square solution.

Use of features in linear regression –

When using data the dependant variables may be affected by more than one independent variable, to make use of this, each independent variable is given weight calculated from the dependant variable and will then be used for prediction

Overfitting and underfitting -

When finding the best fit for the data the aim is to reduce error but there is a point in which the fit is the best it can be. While trying to attain this best fit line underfitting usually occurs, where there is too much error as the line doesn't fit the data. After finding best fit any increased reduction in error will cause overfitting, where the line tries to fit the data so specifically that it is unusable for making predictions.

Intuition of the weight penalty term –

The penalty term/s is used for ridge regression in an attempt to penalise overfitting.

Objective function of ridge regression –

The objective function for ridge regression is the RMSE (Root Mean Squared Error), used to see how closely the functions predictions are to the actual values, this is used on data with all dependant variables already collected.

Ridge Regression Implementation:

When running the code initially, the main function is called. Firstly, the regularization factors are stored in an array and then the following functions are called:

Input, trainassess, ridge_regression, rid_reg_plot and lastly points_plot.

All the functions are discussed below, Input and trainassess are only called once while the rest are called equal to the amount of regularization factors.

```

57 reg_fact_train = np.array([0.000001, 0.0001, 0.01, 0.1]) #10-6, 10-4, 10-2, 10-1
58 (d_tr, d_pl) = input('regression_train_assignment2019.csv', 'regression_plotting_assignment2019.csv')
59 i = 0
60 (tr_x, tr_y, tr_f, pl_x, pl_f) = trainassess(d_tr, d_pl)
61 #(x_tr, y_tr, tr_a) = trainassess(d_tr)
62 for regs in reg_fact_train:
63     (par) = ridge_regression(tr_f, tr_y, regs)
64     (pl_y) = rid_reg_plot(pl_f, par)
65     points_plot(tr_x, tr_y, pl_x, pl_y, regs)
66     i+=1
67
68
69

```

Both Input and trainassess are used for importing data, initially all training and plotting data are placed within Data Frames before being split up into their X values, Features and Y values for training.

```

1  import matplotlib.pyplot as plt
2  import pandas as pd
3  import numpy as np
4
5  def input(train1, test1): #Takes csv files and imports them as dataframes
6      #data_train = pd.DataFrame.from_csv(train1)
7      #data_plot = pd.DataFrame.from_csv(test1)
8      data_train = pd.read_csv(train1)
9      data_plot = pd.read_csv(test1, index_col=None)
10     return data_train, data_plot
11
12
13 def trainassess (d_train, d_plot): #Importing the x,y and features from training. Importing x and Features from plotting
14     x_train = d_train['x'].values
15     y_train = d_train['y'].values
16     f_train = d_train.iloc[:, 3:15].values
17     x_plot = d_plot['x'].values
18     f_plot = d_plot.iloc[:, 2:14].values
19     #print(f_plot)
20     return x_train, y_train, f_train, x_plot, f_plot
21

```

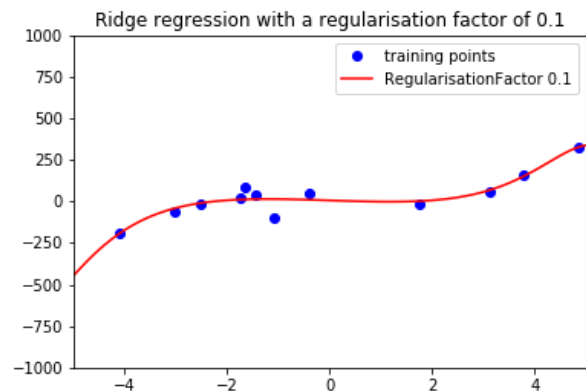
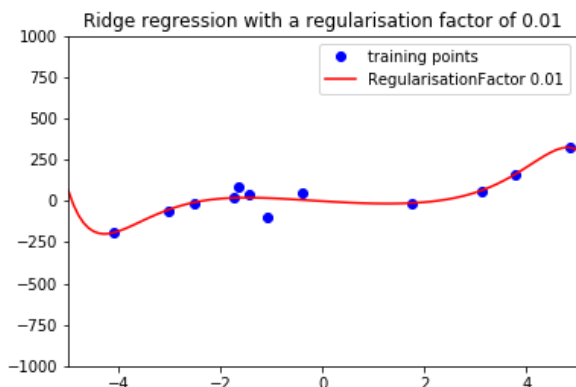
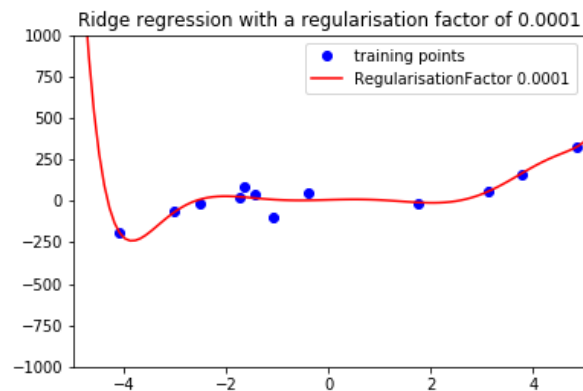
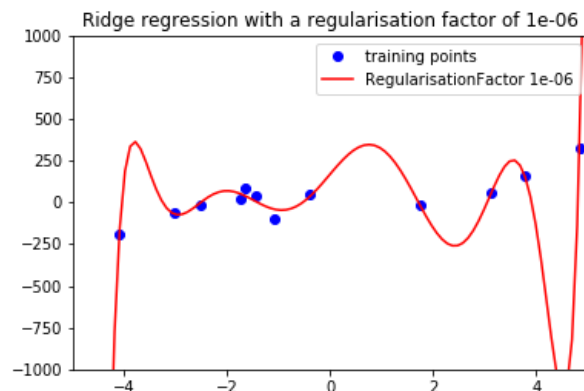
After this the loop starts, initially calling the Ridge Regression equation that calculates the predicted weights each feature must have to give an accurate Y value in this case, this value is changed dependant on the regularization factor. These weights are then used in the next function with the plotting (test) features to predict the y-values for the test data. Finally, the x and y values for both datasets are plotted within the same plot

```

23 def ridge_regression(features_train, y_train, regularisationFactor): #Completes Ridge Regression equation
24
25     parameters = features_train.transpose().dot(features_train)
26     parameters = parameters + (regularisationFactor * (np.identity(parameters.shape[1])))
27     parameters2 = features_train.transpose().dot(y_train)
28     parameters = np.linalg.solve(parameters, parameters2)
29     return parameters
30
31 def rid_reg_plot(f_plot, parameters): #Calculate the predicted Y values by dot producting features with weights
32     i = 0
33     y_calc_plot = np.zeros(len(f_plot))
34     while i < len(f_plot):
35         curr_feat = f_plot[i, :]
36         y_calc_plot[i] = curr_feat.dot(parameters)
37         i+=1
38     return y_calc_plot
39
40 def points_plot(x_train, y_train, x_plot, y_calc_plot, regfact): #Plotting of the training and plot data
41     plt.figure()
42     #plt.clf()
43     plt.plot(x_train, y_train, 'bo')
44     plt.plot(x_plot, y_calc_plot, 'r')
45     #plt.plot(train[0], train[0], 'go')
46     plt.legend(('training points', 'RegularisationFactor ' + str(regfact)))
47     plt.title('Ridge regression with a regularisation factor of ' + str(regfact))
48     plt.xlim((-5, 5))
49     plt.ylim((-1000, 1000))
50     #plt.hold(True)
51     plt.savefig('RidgeRegPlot' + str(regfact) + '.png')
52     plt.show()
53

```

The Plotting above gives the following data for each iteration of the regularization loop:



Ridge Regression Evaluation:

To evaluate the ridge regression a larger list of regularization factors is used, this time the functions called are:

Input, randomtrain, traineval, ridge_regression, eval_regression and lastly plotting_rmse.

Input, randomtrain, traineval and plotting_rmse are called once while the rest are called equal to the amount of regularization factors.

```

62 reg_fact_eval = np.array([0.000001, 0.0001, 0.01, 0.1, 1, 10, 100, 1000])
63 (d_tr, d_pl) = input('regression_train_assignment2019.csv', 'regression_plotting_assignment2019.csv')
64 train_rmse = np.zeros(len(reg_fact_eval))
65 test_rmse = np.zeros(len(reg_fact_eval))
66 (d_tr_rand, d_pl_rand) = randomtrain(d_tr)
67 (tr_x, tr_y, tr_f, pl_x, pl_y, pl_f) = traineval(d_tr_rand, d_pl_rand)
68
69 j = 0
70 for regs in reg_fact_eval:
71     (par) = ridge_regression(tr_f, tr_y, regs)
72     (train_rmse[j]) = eval_regression(par, tr_f, tr_y)
73     (test_rmse[j]) = eval_regression(par, pl_f, pl_y)
74     j+=1
75
76 plotting_rmse(reg_fact_eval, train_rmse, reg_fact_eval, test_rmse)

```

Similar to in the previous implementation, both files are imported but this time only the training data is used, because of its Y values. The data is firstly randomised every iteration and then split in a 70:30 split within the randomtrain function before being outputted back to the main function before being split up within the traineval function, taking x, y and features

```

1  import matplotlib.pyplot as plt
2  import pandas as pd
3  import numpy as np
4
5  def input(train1, test1): #Takes csv files and imports them as dataframes
6      #data_train = pd.DataFrame.from_csv(train1)
7      #data_plot = pd.DataFrame.from_csv(test1)
8      data_train = pd.read_csv(train1)
9      data_plot = pd.read_csv(test1, index_col=None)
10     return data_train, data_plot
11
12 def randomtrain(input_train): #Randomises which part of the train file will be used for training or plotting
13     length = len(input_train)
14     #np.random.shuffle(input_train[0])
15     output = input_train.sample(frac=1)
16     seventy = int(round(0.7*(length)))
17     #print(seventy)
18     output_train = output[0:seventy]
19     output_plot = output[seventy:length]
20     return output_train, output_plot
21
22 def traineval(d_tr, d_pl): #Importing the correct data (features) from the train file specifically
23     x_train = d_tr['x'].values
24     y_train = d_tr['y'].values
25     f_train = d_tr.iloc[:, 3:15].values
26     x_plot = d_pl['x'].values
27     y_plot = d_pl['y'].values
28     f_plot = d_pl.iloc[:, 3:15].values
29     return x_train, y_train, f_train, x_plot, y_plot, f_plot
30

```

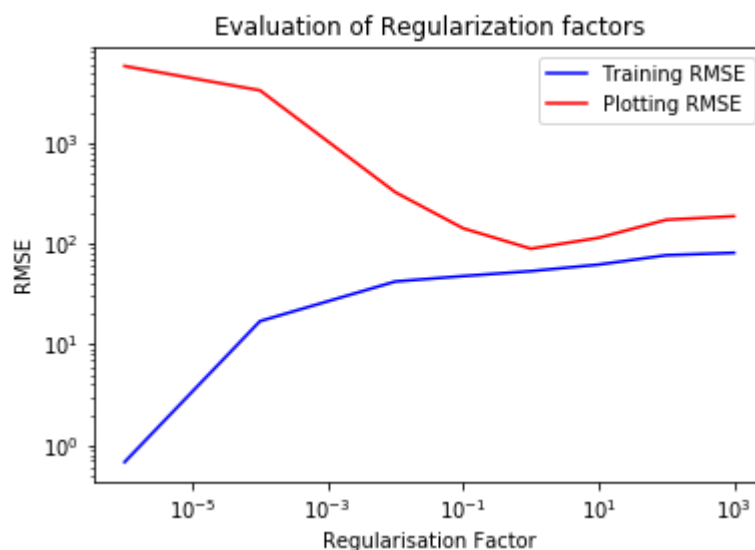
The ridge_regression function is the exact same as during implementation, the new function eval_regression is called twice both for the training and the test data. Firstly, the predicted y-value for the entire dataset is calculated. Then, using this new array, the RMSE (Root Mean Square Error) is found.

```

31 def ridge_regression(features_train, y_train, regularisationFactor): #Computes the weights for the training features
32
33     parameters = features_train.transpose().dot(features_train)
34     parameters = parameters + (regularisationFactor * (np.identity(parameters.shape[1])))
35     parameters2 = features_train.transpose().dot(y_train)
36     parameters = np.linalg.solve(parameters, parameters2)
37     return parameters
38
39 def eval_regression(parameters, features, y): #Calculates the y hat value and then compares it to y to calculate the rmse
40     i = 0
41     y_comp = np.zeros(len(features))
42     while i < len(features):
43         curr_feat = features[i, :]
44         y_comp[i] = curr_feat.dot(parameters)
45         i+=1
46     sq_er = np.square(y - y_comp)
47     mean_sq_er = np.mean(sq_er)
48     rmse = np.sqrt(mean_sq_er)
49     return rmse
50
51 def plotting_rmse(reg_fact_train, train_rmse, reg_fact_eval, test_rmse): #Plots RMSE of both training and testing data
52     plt.figure()
53     plt.loglog(reg_fact_train, train_rmse, 'b')
54     plt.loglog(reg_fact_eval, test_rmse, 'r')
55     plt.title('Evaluation of Regularization factors')
56     plt.legend(('Training RMSE', 'Plotting RMSE'))
57     plt.xlabel('Regularisation Factor')
58     plt.ylabel('RMSE')
59     plt.savefig('RMSEdata.png')
60     plt.show()
61

```

The plotted RMSEs are shown below:



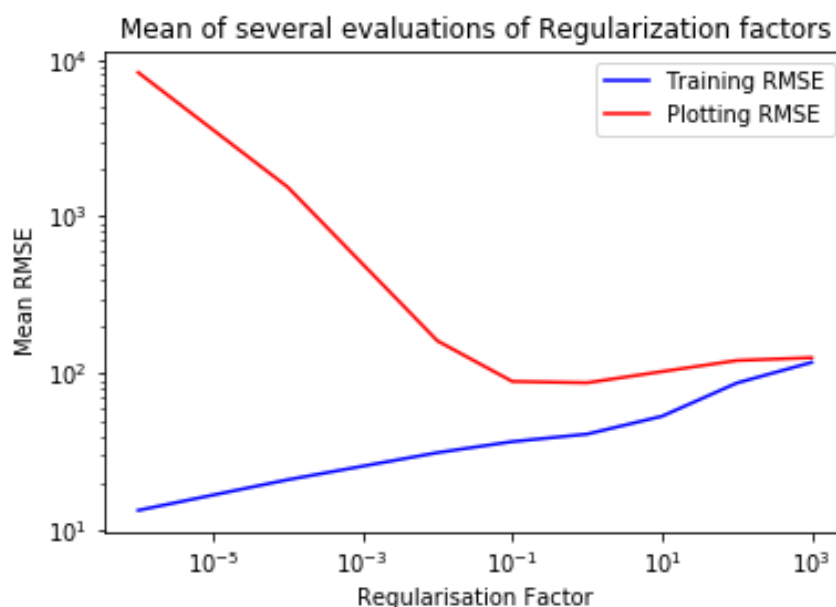
The RMSE's can change quite a lot each time and so a mean of several RMSEs gives a much more accurate display, the only extra code added is the ten loops (rmse_num) that are iterated and the mean values calculated

```

53
54 reg_fact_eval = np.array([0.000001, 0.0001, 0.01, 0.1, 1, 10, 100, 1000])
55 rmse_num = 10
56 train_rmse = np.zeros(len(reg_fact_eval))
57 test_rmse = np.zeros(len(reg_fact_eval))
58 #tr_rmse_mean = np.zeros((rmse_num, len(reg_fact_eval)))
59 tr_rmse_mean = np.zeros(len(reg_fact_eval))
60 #te_rmse_mean = np.zeros((rmse_num, len(reg_fact_eval)))
61 te_rmse_mean = np.zeros(len(reg_fact_eval))
62 i=0
63 while i < rmse_num:
64     (d_tr, d_pl) = input('regression_train_assignment2019.csv', 'regression_plotting_assignment2019.csv')
65     (d_tr, d_pl) = randomtrain(d_tr)
66     (tr_x, tr_y, tr_f, pl_x, pl_y, pl_f) = traineval(d_tr, d_pl)
67     j = 0
68     for regs in reg_fact_eval:
69         (par) = ridge_regression(tr_f, tr_y, regs)
70         (train_rmse[j]) = eval_regression(par, tr_f, tr_y)
71         (test_rmse[j]) = eval_regression(par, pl_f, pl_y)
72         tr_rmse_mean[j] += train_rmse[j]
73         te_rmse_mean[j] += test_rmse[j]
74         j+=1
75     i+=1
76 tr_mean = tr_rmse_mean / rmse_num
77 te_mean = te_rmse_mean / rmse_num
78 #tr_mean = np.mean
79 plotting_rmse(reg_fact_eval, tr_mean, reg_fact_eval, te_mean)

```

The plotting shows that given the data 10 to the power of -1 is the best Regularization factor:



Task 2**K-Means Clustering Description:****Centroids –**

Centroids are the cluster centre points in which data can be assigned and so the number of centroids is equal to the number of clusters needed.

Euclidean distance –

This is the distance between centroids and each data point, taking the square distance of all features from the centroids. The square is taken to make sure that any distances that are lower than the centroid in any dimension do not reduce the Euclidean distance unnecessarily. This value is also used to calculate errors.

Assignment step –

The data points are all assigned to the closest centroid, calculated via the Euclidean distance.

Update step –

Once all data points have been assigned to their closest centroids the mean of the of their values are taken, this will be the new centroid, iteration follows until the new centroids have not changed from the previous and a convergence is reached, the clusters are then which data points are closest to each of the final centroids.

Objective function –

The objective function is the value we are attempting to reduce and use to evaluate how well the clustering fits the data. For this the RSS for each iteration of values of k can show that the errors is decreasing and the clustering is not getting further away, if the objective is to find out how many clusters are necessary and by extension what the best value of k is, then the final RSS can be found for each k value and compared. If as the number of clusters increases the aggregate distance decreases steeply then at the point the gradient becomes shallow or where the “mountain ends” and the “rubbles’ begins” (Zhong, M. 2019)

K-Means Clustering Implementation:

The code loops for the two k-values, each time importing the csv file, this first set of code was used to ensure no issues occurred when using more than two features so stem_data only stores stem_length and stem_diameter. Later code will use all features. The kmeans function is then called which will itself call several functions, the only other functions are for plotting both the clusters in a two dimensional space and to show the objective.

```

143
144 k = [3, 4]
145 i = 0
146 while i < len(k):
147     #data_train = np.loadtxt(open("CMP3744M_ADM_Assignment 1_Task2 - dataset - plants.csv", "rb"), delimiter=",")
148     Cluster_data = pd.read_csv('CMP3744M_ADM_Assignment 1_Task2 - dataset - plants.csv', index_col=None)
149     stem_data = Cluster_data.iloc[:,0:2].values
150
151     (cent, clu_as, iterat, obj) = kmeans(stem_data, k[i])
152     #print(clu_as)
153     #print(cent)
154     x_1 = Cluster_data.iloc[:,0].values
155     y_1 = Cluster_data.iloc[:,1].values
156     x_2 = Cluster_data.iloc[:,2].values
157     y_2 = Cluster_data.iloc[:,3].values
158     #x_train = data_train[:,0]
159     #y_train = data_train[:,1]
160     iteration = np.array(range(1,(iterat+1)))
161     #print(iteration)
162     #k = KMeans(n_clusters=3).fit_predict(stem_data)
163     #plt.scatter(x_train1, x_train2, c=k)
164     centx = cent[:,0]
165     centy = cent[:,1]
166     #print(stem_data[:,0].shape)
167     plot_clustering(x_1,y_1, clu_as, centx, centy, k[i])
168     plot_iter_obj(iteration, obj, k[i])
169     i+=1

```

Within the kmeans function a number of values are set including which cluster the variables have been assigned to, a comparison cluster array, the total of the Euclidean distances for each data point and the sum of all Euclidean distances.

This function calls several functions within itself, including:

 Initialise_centroids and compute_euclidean_distance

The initialise_centroids is called at the start and should only be called once if the centroids are allocated correctly, compute_euclidean_distance is called every time the distance from the centroids is needed, they will be discussed in more detail below.


```

47
48 def kmeans(dataset, k): # k=3 or 4
49     chk_centroids_zero = True
50     #objective = 0.0
51     while chk_centroids_zero == True:
52         clust_assi = np.ones(dataset[:,0].size) #Current cluster assigned to the dataset
53         clust_comp = np.zeros(dataset[:,0].size) #Previous cluster assigned to the dataset, used for comparison
54         objective = np.array([])
55         objective_sum = np.array([]) #Used to hold the objective function
56         cluster_assigned = np.empty(dataset[:,0].size) #Once completed clusters will be assigned here
57         euc_xy_tot = np.zeros([k, 2]) #The total euclidian distance for x and y
58         euc_clus_tot = np.zeros(k)
59         mean_tot = np.zeros([k, 2])
60         loops = 0
61         #while chk_centroids_zero == True:
62         centroids = initialise_centroids(dataset, k)
63         #while
64         #print(centroids)
65         while np.array_equal(clust_assi, clust_comp) != True:
66             objective_sum = np.array([]) #Reset Calculations of Euclidean error each loop
67             euc_xy_tot = np.zeros([k, 2])
68             euc_clus_tot = np.zeros(k)
69             #print(clust_assi, clust_comp)
70             #clust_comp = np.zeros(dataset[:,0].size)
71             clust_comp[:] = clust_assi[:]
72             #print(clust_comp)
73             i = 0
74             while i < dataset[:,0].size : #Loops through all 300 values
75
76                 dist_data = dataset[i, :] # takes the ith x and y values
77                 #print(dist_data)
78                 #calculates the euc dist for each datapoint per centroid
79                 dist = compute_euclidean_distance(dist_data, centroids)
80                 #print (dist)
81                 mindist = np.argmin(dist)
82                 shortest_dist = dist[mindist]
83                 #print(shortest_dist)
84                 #print(mindist)
85                 objective_sum = np.append(objective_sum, shortest_dist)
86                 euc_xy_tot[mindist, 0] += dataset[i, 0]
87                 euc_xy_tot[mindist, 1] += dataset[i, 1]
88                 clust_assi[i] = mindist
89                 #print(clust_assi[i])
90                 euc_clus_tot[mindist] += 1
91                 i+=1
92             objective = np.append(objective, np.sum(objective_sum))
93
94
95             if 0 not in euc_clus_tot: #Used to make sure the centroids have data points within them
96                 chk_centroids_zero = False
97                 #print(objective)
98             j=0
99             while j < k :
100                 if euc_xy_tot[j, 0] != 0 and euc_xy_tot[j, 1] != 0:
101                     mean_tot[j, 0] = euc_xy_tot[j, 0] / euc_clus_tot[j]
102                     mean_tot[j, 1] = euc_xy_tot[j, 1] / euc_clus_tot[j]
103                 j+=1
104             #print(mean_tot)
105             centroids[:] = mean_tot[:]
106             #print(centroids)
107             #print(euc_xy_tot)
108             #print(clust_assi, clust_comp)
109
110             loops+=1
111             #print(objective)
112             cluster_assigned = clust_assi[:]
113         return centroids, cluster_assigned, loops, objective
114

```

Within the initialise_centroids function random float values are chosen from between the upper and lower bounds of the dataset, this is to reduce the chance the centroids will be so far outside the dataset that no data points will be assigned to them, leading to the eventual centroid being at 0,0 on a two dimensional graph. The compute_euclidean_distance function is called for each data point and compares it to all centroids and their features.

```

1 import matplotlib.pyplot as plt
2 import pandas as pd
3 import numpy as np
4 import random as rd
5
6
7 def initialise_centroids(dataset, k): #k=3 or 4, called once at start of code
8     centroids = np.empty([k, 2]) #Create an empty array
9     i = 0
10    kmax = np.zeros([k])
11    #dtmin = int(np.min(dataset))
12    #dtmax = int(np.max(dataset))
13    dtmin = np.min(dataset)
14    dtmax = np.max(dataset)
15
16    #dtminx = np.argmin(dataset[:,0])
17    #dtminx2 = int(dataset[dtminx, 0])
18    #dtmaxx = np.argmax(dataset[:,0])
19    #dtmaxx2 = int(dataset[dtmaxx, 0])
20
21    #dtminy = np.argmin(dataset[:,1])
22    #dtminy2 = int(dataset[dtminy, 1])
23    #dtmaxy = np.argmax(dataset[:,1])
24    #dtmaxy2 = int(dataset[dtmaxy, 1])
25
26    while i < k: #going through array until all centroids have a random x and y value
27        #centroids[i,0] = rd.randint(dtmin,dtmax)
28        #centroids[i,1] = rd.randint(dtmin,dtmax)
29        centroids[i,0] = rd.uniform(dtmin,dtmax)
30        centroids[i,1] = rd.uniform(dtmin,dtmax)
31        i+=1
32    return centroids
33
34
35 def compute_euclidean_distance(dtst, centr):
36     distance = np.empty([centr[:,0].size])
37     #print(test)
38     i = 0
39     while i < centr[:,0].size :
40         #Euclidian distance = the sum of data point values - centroid values then squared
41         #distance[i] = np.sqrt(np.square(dtst[0] - centr[i,0]) + np.square(dtst[1] - centr[i,1]))
42         distance[i] = np.sqrt(np.square(dtst[0] - centr[i,0]) + np.square(dtst[1] - centr[i,1]))
43         i+=1
44     #distance = np.sum
45     return distance
46

```

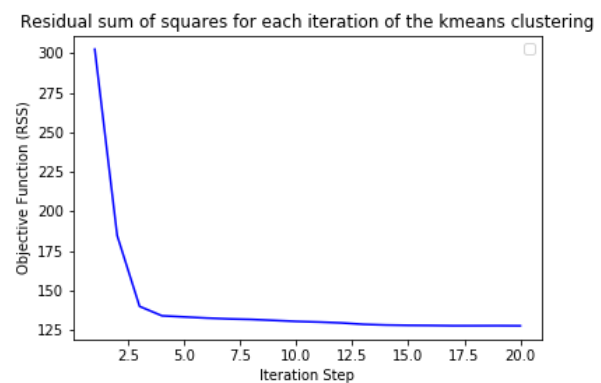
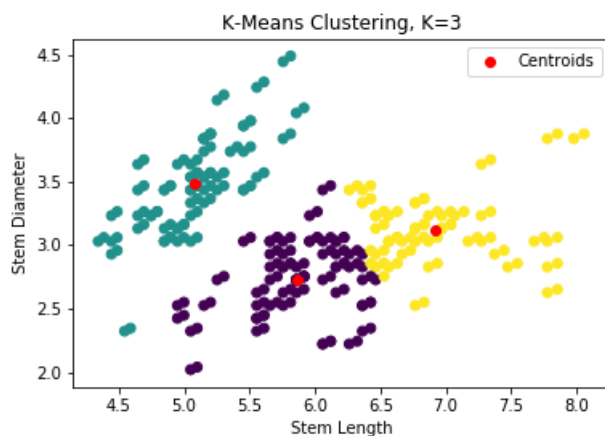
The plotting functions called after kmeans completes display the data below

```

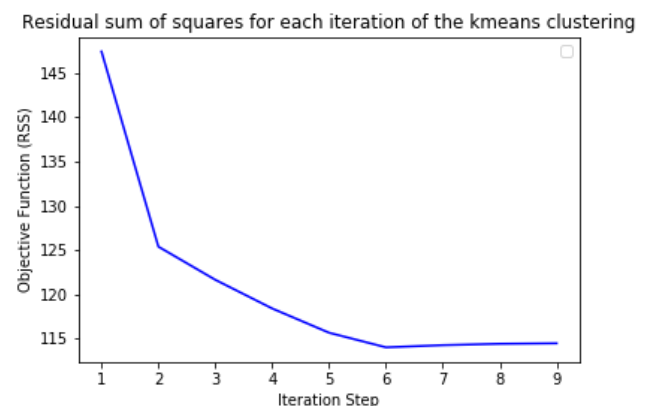
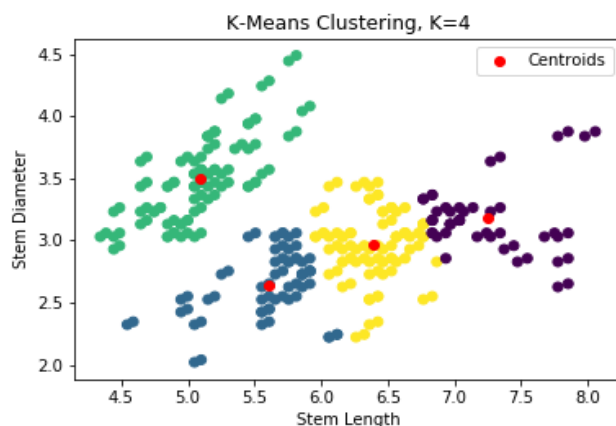
114
115 def plot_clustering(x1,y1, clus, centx, centy, k):
116     plt.figure()
117     plt.scatter(x1, y1, c=clus)
118     plt.scatter(centx, centy, c='r', label='Centroids')
119     plt.legend()
120     plt.title('K-Means Clustering, K=' + str(k))
121     plt.xlabel('Stem Length')
122     plt.ylabel('Stem Diameter')
123     plt.savefig('TwoFeatureCluster' + str(k) + '.png')
124     plt.show()
125
126 def plot_iter_obj(iterations, obj_func, k): # inputs iterations, the objective function (euclidean)
127     plt.figure()
128     plt.plot(iterations, obj_func, c='b')
129     plt.legend()
130     plt.title('Residual sum of squares for each iteration of the kmeans clustering')
131     plt.xlabel('Iteration Step')
132     plt.ylabel('Objective Function (RSS)') #Sum of the sum of the euclidian distances per iteration
133     plt.savefig('TwoFeatureRSS' + str(k) + '.png')
134     plt.show()
135

```

For K = 3 the plots show:



For K = 4 the plots show:



To allow for full use of the kmeans all four features must be used to group the data, small changes were made to all functions to be resilient to as many features as necessary, the entirety of the code will be shown below

```

1 import matplotlib.pyplot as plt
2 import pandas as pd
3 import numpy as np
4 import random as rd
5
6
7 def initialise_centroids(dataset, k): #k=3 or 4, called once at start of code
8     num_of_feat = dataset[0,:].size
9     centroids = np.empty([k, num_of_feat]) #Create an empty array
10    i = 0
11    j = 0
12    kmax = np.zeros([k])
13    dtmin = np.min(dataset)
14    dtmax = np.max(dataset)
15    while j < num_of_feat:
16        while i < k: #going through array until all centroids have a random x and y value
17            centroids[i,j] = rd.uniform(dtmin,dtmax)
18            i+=1
19        i=0
20        j+=1
21    return centroids
22
23
24 def compute_euclidean_distance(dtst, centr):
25     #print(dtst)
26     num_of_feat = dtst.size
27     distance = np.empty([centr[:,0].size])
28     #print(test)
29     i = 0
30     j = 0
31     while i < centr[:,0].size:
32         while j < num_of_feat:
33             #Euclidian distance = the sum of data point values - centroid values then squared
34             #distance[i] = np.sqrt(np.square(dtst[0] - centr[i,0]) + np.square(dtst[1] - centr[i,1]))
35             distance[i] += np.square(dtst[j] - centr[i,j])
36             #print(distance[i])
37             j+=1
38         distance[i] = np.sqrt(distance[i])
39         j=0
40         i+=1
41     #print(distance)
42     return distance
43

```

```

45 def kmeans(dataset, k): # k=3 or 4
46     chk_centroids_zero = True
47     num_of_feat = dataset[0,:].size
48     #objective = 0.0
49     while chk_centroids_zero == True:
50         clust_assi = np.ones(dataset[:,0].size) #Current cluster assigned to the dataset
51         clust_comp = np.zeros(dataset[:,0].size) #Previous cluster assigned to the dataset, used for comparison
52         objective = np.array([]) #Used to hold the objective function
53         objective_sum = np.array([]) #Used to hold the objective function
54         cluster_assigned = np.empty(dataset[:,0].size) #Once completed clusters will be assigned here
55         euc_xy_tot = np.zeros([k, num_of_feat]) #The total euclidian distance for all features
56         euc_clus_tot = np.zeros(k)
57         mean_tot = np.zeros([k, num_of_feat])
58         loops = 0
59         #while chk_centroids_zero == True:
60         centroids = initialise_centroids(dataset, k)
61         #while
62         #print(centroids)
63         while np.array_equal(clust_assi, clust_comp) != True:
64             objective_sum = np.array([]) #Reset Calculations of Euclidean error each Loop
65             euc_xy_tot = np.zeros([k, num_of_feat])
66             euc_clus_tot = np.zeros(k)
67             #print(clust_assi, clust_comp)
68             #clust_comp = np.zeros(dataset[:,0].size)
69             clust_comp[:] = clust_assi[:]
70             #print(clust_comp)
71             i = 0
72             while i < dataset[:,0].size : #Loops through all 300 values
73
74                 dist_data = dataset[i, :] # takes the ith features
75                 #print(dist_data)
76                 dist = compute_euclidean_distance(dist_data, centroids) #calculates the euc dist for each datapoint per centi
77                 #print (dist)
78                 mindist = np.argmin(dist)
79                 shortest_dist = dist[mindist]
80                 #print(shortest_dist)
81                 #print(mindist)
82                 objective_sum = np.append(objective_sum, shortest_dist)
83                 euc_xy_tot[mindist, 0] += dataset[i, 0]
84                 euc_xy_tot[mindist, 1] += dataset[i, 1]
85                 clust_assi[i] = mindist
86                 #print(clust_assi[i])
87                 euc_clus_tot[mindist] += 1
88                 i+=1
89             objective = np.append(objective, np.sum(objective_sum))
90             #print(objective)
91
92             if 0 not in euc_clus_tot: #Used to make sure the centroids have data points within them
93                 chk_centroids_zero = False
94                 #print(objective)
95             j=0
96             while j < k :
97                 if euc_xy_tot[j, 0] != 0 and euc_xy_tot[j, 1] != 0:
98                     mean_tot[j, 0] = euc_xy_tot[j, 0] / euc_clus_tot[j]
99                     mean_tot[j, 1] = euc_xy_tot[j, 1] / euc_clus_tot[j]
100                 j+=1
101             #print(mean_tot)
102             centroids[:] = mean_tot[:]
103             #print(centroids)
104             #print(euc_xy_tot)
105             #print(clust_assi, clust_comp)
106
107             loops+=1
108             #print(objective)
109             cluster_assigned = clust_assi[:]
110         return centroids, cluster_assigned, loops, objective
111

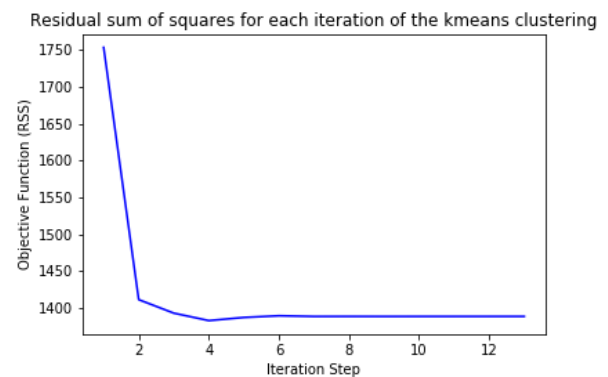
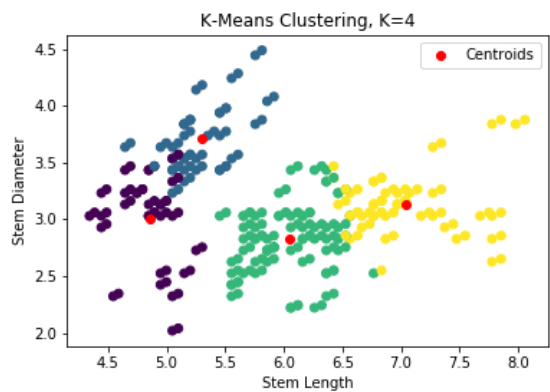
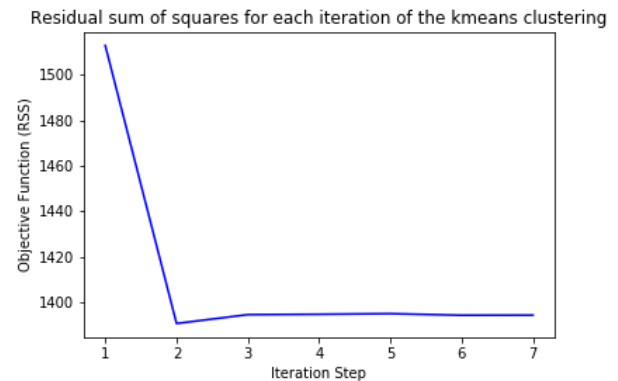
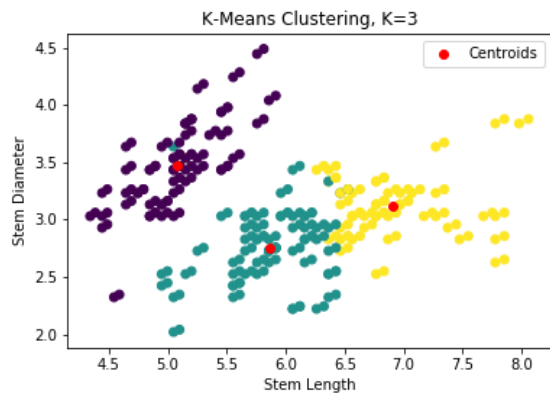
```

```

111
112 def plot_clustering(x1,y1, clus, centx, centy, k):
113     plt.figure()
114     plt.scatter(x1, y1, c=clus)
115     plt.scatter(centx, centy, c='r', label='Centroids')
116     plt.legend()
117     plt.title('K-Means Clustering, K=' + str(k))
118     plt.xlabel('Stem Length')
119     plt.ylabel('Stem Diameter')
120     plt.savefig('FullCluster' + str(k) + '.png')
121     plt.show()
122
123 def plot_iter_obj(iterations, obj_func, k): # inputs iterations, the objective function (euclidean)
124     plt.figure()
125     plt.plot(iterations, obj_func, c='b')
126     #plt.legend()
127     plt.title('Residual sum of squares for each iteration of the kmeans clustering')
128     plt.xlabel('Iteration Step')
129     plt.ylabel('Objective Function (RSS)') #Sum of the sum of the euclidian distances per iteration
130     plt.savefig('FullRSS' + str(k) + '.png')
131     plt.show()
132
133 k = [3, 4]
134 i = 0
135 while i < len(k):
136     #data_train = np.loadtxt(open("CMP3744M_ADM_Assignment 1_Task2 - dataset - plants.csv", "rb"), delimiter=",")
137     Cluster_data = pd.read_csv('CMP3744M_ADM_Assignment 1_Task2 - dataset - plants.csv', index_col=None)
138     all_data = Cluster_data.iloc[:,0:4].values
139
140     (cent, clu_as, iterat, obj) = kmeans(all_data, k[i])
141     #print(clu_as)
142     #print(cent)
143     x_1 = Cluster_data.iloc[:,0].values
144     y_1 = Cluster_data.iloc[:,1].values
145     x_2 = Cluster_data.iloc[:,2].values
146     y_2 = Cluster_data.iloc[:,3].values
147     #x_train = data_train[:,0]
148     #y_train = data_train[:,1]
149     iteration = np.array(range(1,(iterat+1)))
150     #print(iteration)
151     #k = KMeans(n_clusters=3).fit_predict(stem_data)
152     #plt.scatter(x_train1, x_train2, c=k)
153     centx = cent[:,0]
154     centy = cent[:,1]
155     #print(stem_data[:,0].shape)
156     plot_clustering(x_1,y_1, clu_as, centx, centy, k[i])
157     plot_iter_obj(iteration, obj, k[i])
158     i+=1

```

It produces slightly different results shown below, there is some overlap with different clusters as the features are create a space that is a higher dimension than that used to plot the graphs:



References:

Bishop, C.M. (2006) *Pattern recognition and machine learning*. Oxford: Springer
 Goodfellow, I., Bengio, Y. and Courville, A. (2016) *Deep Learning*. Cambridge, Massachusetts: The MIT Press
 Han, J., Kamber, M. and Pei, J. (2012) *Data Mining: Concepts and Techniques*. Amsterdam: Morgan Kaufmann
 James, G., Witten, D., Hastie, T. and Tibshirani, R. (2013) *An Introduction to Statistical Learning with Applications in R*. New York, United States: Springer
 Zhong, M. (2019) *Clustering* [lecture]. Algorithms and Data Mining CMP3744M-1819, University of Lincoln, 21 February. Available from: https://blackboard.lincoln.ac.uk/bbcswebdav/pid-2201579-dt-content-rid-3878326_2/courses/CMP-ADM-1819/ADM_W5_Lec_Clustering_V2.pdf [accessed 17 March 2019]