

1 Polynomial Regression (Programming)

```
1 # Machine Learning HW2 Poly Regression
2
3 import matplotlib.pyplot as plt
4 import numpy as np
5
6 # Step 1
7 # Parse the file and return 2 numpy arrays
8
9
10 def load_data_set(filename):
11     arr = np.loadtxt(filename)
12     x, y = np.split(arr, [-1], axis=1)
13     plt.plot(x, y, '.')
14     plt.show()
15     return x, y
16
17 # Find theta using the normal equation
18
19
20 def normal_equation(x, y):
21     x_t = np.transpose(x)
22     theta = np.dot(np.dot(np.linalg.inv(np.dot(x_t, x)), x_t), y)
23     return theta, []
24
25 # Step 2:
26 # Given a n by 1 dimensional array return an n by num_dimension array
27 # consisting of [1, x, x^2, ...] in each row
28 # x: input array with size n
29 # degree: degree number, an int
30 # result: polynomial basis based reformulation of x
31
32
33 def increase_poly_order(x, degree):
34     result = np.array([list(np.power(x.flatten(), i))
35                        for i in range(degree+1)]).T
36     # normalize
37     # result = result / result.max(axis=0)
38     return result
39
40 # split the data into train and test examples by the train_proportion
41 # i.e. if train_proportion = 0.8 then 80% of the examples are training and 20%
42 # are testing
43
44
45 def train_test_split(x, y, train_proportion):
46     # your code
47     num_train = int(x.shape[0] * train_proportion)
48     x_train, x_test = x[:num_train, :], x[num_train:, :]
49     y_train, y_test = y[:num_train, :], y[num_train:, :]
50     return x_train, x_test, y_train, y_test
51
52 # Find theta using the gradient descent
53
54
55 def solve_regression(x, y, num_iterations=100, learning_rate=0.002):
56     # your GD code from HW1 or better version
57     num_features = x.shape[1]
58     thetas = []
59     theta = np.array([[0] for i in range(num_features)])
60     for _ in range(num_iterations):
61         gradient_sum = np.array([[0.0] for i in range(num_features)])
62         for row, label in zip(x, y):
63             gradient_sum = gradient_sum - learning_rate *
64                 row[:, np.newaxis] * (np.dot(row, theta) - label)
65         theta = theta + gradient_sum
66         thetas.append(theta)
67     return theta, thetas
68
```

```

69 # Given an array of y and y_predict return loss
70 # y: an array of size n
71 # y_predict: an array of size n
72 # loss: a single float
73
74
75 def get_loss(y, y_predict):
76     diff = y - y_predict
77     loss = np.dot(diff.T, diff) / len(y)
78     return loss[0][0]
79
80 # Given an array of x and theta predict y
81 # x: an array with size n x d
82 # theta: np array including parameters
83 # y_predict: prediction labels, an array with size n
84
85
86 def predict(x, theta):
87     # your code
88     y_predict = x.dot(theta)
89     return y_predict
90
91
92 # Given a list of thetas one per (s)GD epoch
93 # this creates a plot of epoch vs prediction loss (one about train, and another about test)
94 # this figure checks GD optimization traits of the best theta
95 def plot_epoch_losses(x_train, x_test, y_train, y_test, best_thetas, title):
96     # your code
97     epochs = []
98     losses = []
99     epoch_num = 1
100     for theta in best_thetas:
101         losses.append(get_loss(y_train, predict(x_train, theta)))
102         epochs.append(epoch_num)
103         epoch_num += 1
104     fig, ax = plt.subplots()
105     ax.scatter(epochs, losses, label="training")
106     plt.xlabel("epoch")
107     plt.ylabel("loss")
108     plt.title(title)
109     # plt.show()
110
111     epochs = []
112     losses = []
113     epoch_num = 1
114     for theta in best_thetas:
115         losses.append(get_loss(y_test, predict(x_test, theta)))
116         epochs.append(epoch_num)
117         epoch_num += 1
118     ax.scatter(epochs, losses, label="testing")
119     ax.legend()
120     # plt.xlabel("epoch")
121     # plt.ylabel("loss")
122     # plt.title(title + "on testing data")
123     plt.show()
124
125
126 # Given a list of degrees.
127 # For each degree in the list, train a polynomial regression.
128 # Return training loss and validation loss for a polynomial regression of order degree for
129 # each degree in degrees.
130 # Use 60% training data and 20% validation data. Leave the last 20% for testing later.
131 # Input:
132 # x: an array with size n x d
133 # y: an array with size n
134 # degrees: A list of degrees
135 # Output:
136 # training_losses: a list of losses on the training dataset
137 # validation_losses: a list of losses on the validation dataset
138 def get_loss_per_poly_order(x, y, degrees):
139     # your code
140     training_losses = []
141     validation_losses = []
142     for degree in degrees:

```

```

143     augmented_x = increase_poly_order(x, degree)
144     x_train, x_test, y_train, y_test = train_test_split(
145         augmented_x, y, 0.6)
146     x_validation, x_test, y_validation, y_test = train_test_split(
147         x_test, y_test, 0.5)
148     theta, thetas = normal_equation(x_train, y_train)
149     training_losses.append(get_loss(y_train, predict(x_train, theta)))
150     validation_losses.append(
151         get_loss(y_validation, predict(x_validation, theta)))
152     return training_losses, validation_losses
153
154 # Give the parameter theta, best-fit degree, plot the polynomial curve
155
156
157 def best_fit_plot(theta, degree):
158     # your code
159     augmented_x = increase_poly_order(x, degree)
160     y_predict = predict(augmented_x, theta)
161     sorted_y_predict = [y for _, y in sorted(zip(x[:, 0], y_predict))]
162     plt.scatter(x[:, 0], y)
163     plt.plot(sorted(x[:, 0]), sorted_y_predict, 'y')
164     plt.xlabel("X")
165     plt.ylabel("y")
166     plt.title("Scatter Plot of Data")
167     plt.show()
168
169
170 def select_hyperparameter(degrees, x_train, x_test, y_train, y_test):
171     # Part 1: hyperparameter tuning:
172     # Given a set of training examples, split it into train-validation splits
173     # do hyperparameter tune
174     # come up with best model, then report error for best model
175     training_losses, validation_losses = get_loss_per_poly_order(
176         x_train, y_train, degrees)
177     plt.plot(degrees, training_losses, label="training_loss")
178     plt.plot(degrees, validation_losses, label="validation_loss")
179     plt.yscale("log")
180     plt.legend(loc='best')
181     plt.title("poly order vs validation_loss")
182     plt.show()
183
184     # Part 2: testing with the best learned theta
185     # Once the best hyperparameter has been chosen
186     # Train the model using that hyperparameter with all samples in the training
187     # Then use the test data to estimate how well this model generalizes.
188     best_degree = 5 # fill in using best degree from part 2
189     x_train = increase_poly_order(x_train, best_degree)
190     best_theta, best_thetas = normal_equation(x_train, y_train)
191     print(best_theta)
192     best_fit_plot(best_theta, best_degree)
193     x_test = increase_poly_order(x_test, best_degree)
194     test_loss = get_loss(y_test, predict(x_test, best_theta))
195     train_loss = get_loss(y_train, predict(x_train, best_theta))
196     return best_degree, best_theta, train_loss, test_loss
197
198
199 # Given a list of dataset sizes [d_1, d_2, d_3 .. d_k]
200 # Train a polynomial regression with first d_1, d_2, d_3, .. d_k samples
201 # Each time,
202 # return the a list of training and testing losses if we had that number of examples.
203 # We are using 0.5 as the training proportion because it makes the testing_loss more stable
204 # in reality we would use more of the data for training.
205 # Input:
206 # x: an array with size n x d
207 # y: an array with size n
208 # example_num: A list of dataset size
209 # Output:
210 # training_losses: a list of losses on the training dataset
211 # testing_losses: a list of losses on the testing dataset
212 def get_loss_per_tr_num_examples(x, y, example_num, train_proportion):
213     # your code
214     print(x.shape)
215     training_losses = []
216     testing_losses = []

```

```

217     for n in example_num:
218         x_available, y_available = x[:n], y[:n]
219         x_train, x_test, y_train, y_test = train_test_split(
220             x_available, y_available, train_proportion)
221         theta, thetas = normal_equation(x_train, y_train)
222         training_losses.append(get_loss(y_train, predict(x_train, theta)))
223         testing_losses.append(get_loss(y_test, predict(x_test, theta)))
224     return training_losses, testing_losses

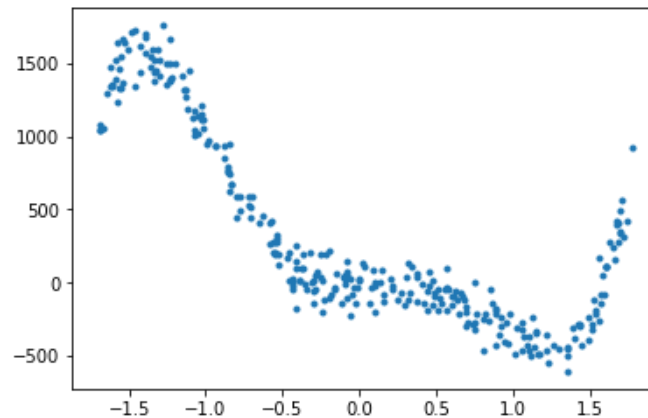
```

Loading Data Set

```

1  x, y = load_data_set("dataPoly.txt")

```

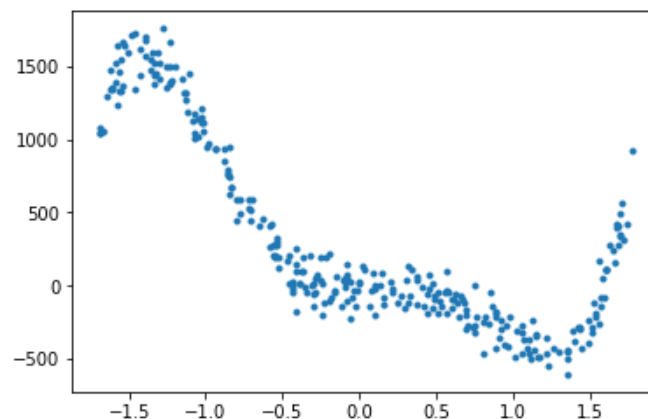


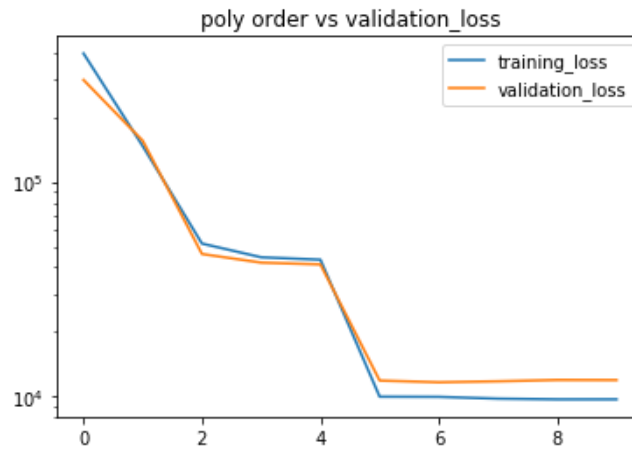
Task1 & Task2

```

1  # select the best polynomial through train-validation-test formulation
2  x, y = load_data_set("dataPoly.txt")
3  x_train, x_test, y_train, y_test = train_test_split(x, y, 0.8)
4  degrees = [i for i in range(10)]
5  best_degree, best_theta, train_loss, test_loss = select_hyperparameter(
6      degrees, x_train, x_test, y_train, y_test)

```

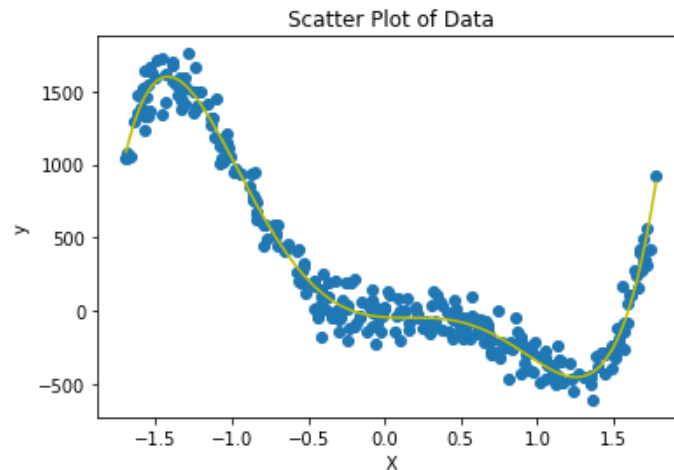




```

1  [[ -42.7987349 ]
2   [ -73.48160041]
3   [ 454.86959156]
4   [-927.86375004]
5   [ -63.57127418]
6   [ 307.7877446  ]]

```



As we can see in the poly order vs validation loss figure, increasing the degree beyond 5 does not decrease the validation loss, which makes sense because y is produced by a quintic (degree 5) polynomial function

With the degree set to 5, the best theta is

$[-4.49472078e + 01, -7.59089828e + 00, 1.06657223e + 01, -4.93874067e + 00, -5.94462091e - 03, 4.94797665e - 02]^T$, which is reasonably close to the data generation function $0.05x^5 - 5x^3 + 10x^2 - 5x - 30$

Task3

```

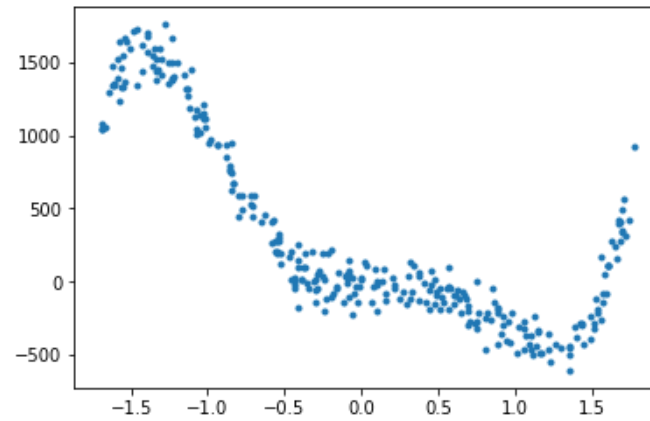
1  def trainGD(degrees, x_train, x_test, y_train, y_test, best_degree = 5):
2      # Part 2: testing with the best learned theta
3      # Once the best hyperparameter has been chosen
4      # Train the model using that hyperparameter with all samples in the training
5      # Then use the test data to estimate how well this model generalizes.
6      x_train = increase_poly_order(x_train, best_degree)
7      best_theta, best_thetas = solve_regression(x_train, y_train, 200, 0.0002)
8      print(best_theta)
9      best_fit_plot(best_theta, best_degree)
10     x_test = increase_poly_order(x_test, best_degree)
11     test_loss = get_loss(y_test, predict(x_test, best_theta))
12     train_loss = get_loss(y_train, predict(x_train, best_theta))
13
14     # Part 3: visual analysis to check GD optimization traits of the best theta
15     plot_epoch_losses(x_train, x_test, y_train, y_test, best_thetas,
16                       "best learned theta - train, test losses vs. GD epoch ")
17     return best_degree, best_theta, train_loss, test_loss
18

```

```

19 # select the best polynomial through train-validation-test formulation
20 x, y = load_data_set("dataPoly.txt")
21 x_train, x_test, y_train, y_test = train_test_split(x, y, 0.8)
22 degrees = [i for i in range(10)]
23 print("degree=5")
24 best_degree, best_theta, train_loss, test_loss = trainGD(
25     degrees, x_train, x_test, y_train, y_test)
26 print("degree=3")
27 best_degree, best_theta, train_loss, test_loss = trainGD(
28     degrees, x_train, x_test, y_train, y_test, 3)

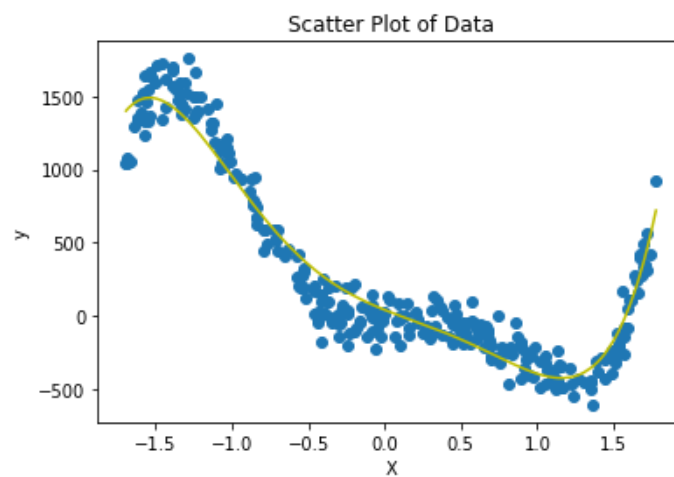
```

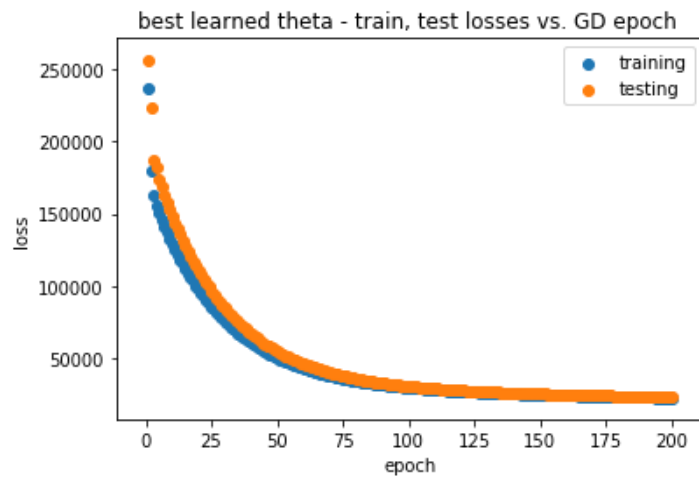


```

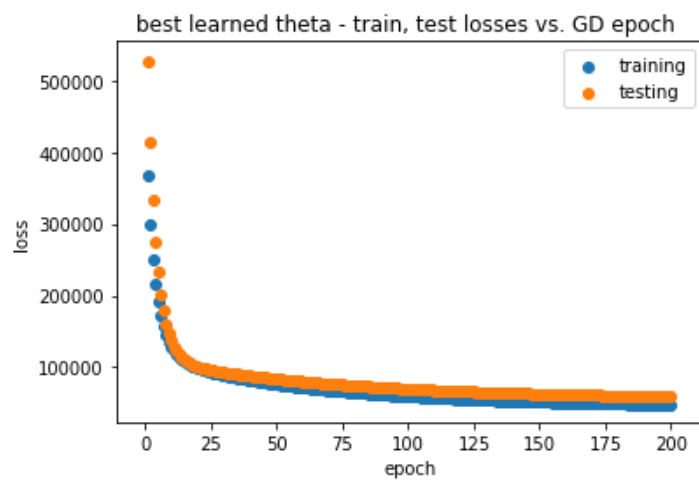
1 degree=5
2 [[ 41.52331804]
3  [-419.28587676]
4  [ 210.20916279]
5  [-413.00549955]
6  [ 28.46261813]
7  [ 157.56251091]]

```





```
1 degree=3
2 [[ -17.20149937]
3  [-531.21580312]
4  [ 335.34485717]
5  [ 12.05204491]]
```



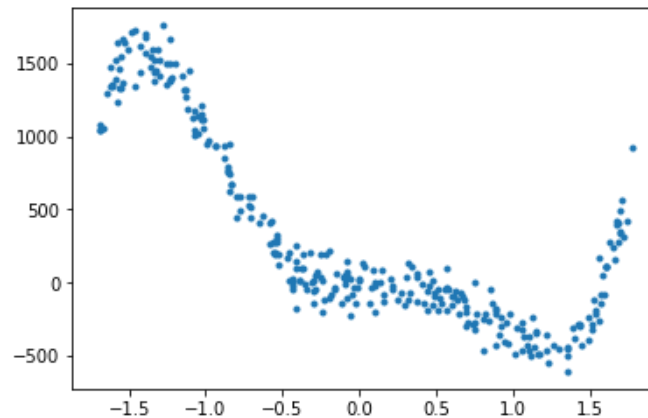
The above figure shows that if we choose too small a degree, the difference between the training error and the testing error will be a little larger because the model isn't complex enough to capture the patterns in the data

Task4

```

1 # Part 4: analyze the effect of revising the size of train data:
2 # Show training error and testing error by varying the number for training samples
3 x, y = load_data_set("dataPoly.txt")
4 x = increase_poly_order(x, 8)
5 example_num = [10*i for i in range(2, 11)] # python list comprehension
6 training_losses, testing_losses = get_loss_per_tr_num_examples(
7     x, y, example_num, 0.5)
8 plt.plot(example_num, training_losses, label="training_loss")
9 plt.plot(example_num, testing_losses, label="testing_losses")
10 plt.yscale("log")
11 plt.legend(loc='best')
12 plt.title("number of examples vs training_loss and testing_loss")
13 plt.show()

```



```

1 (300, 9)

```



In the figure above, we can see testing loss reaches its minimum at $n = 40$. After that, more training examples didn't help improving the performance of the model. I suspect this is because we set the degree to 8, which is higher than that of the data generation function, making the model prone to overfitting.

2 Ridge Regression (programming and QA)

2.1 QA

1.1

$$\frac{\partial J(\beta)}{\partial \beta} = \frac{y^T - y^T X \beta - \beta^T X^T y + \beta^T X^T X \beta + \lambda \beta^T \beta}{\partial \beta} = \frac{\beta^T X^T X \beta - 2\beta^T X^T y + \beta^T (\lambda I) \beta}{\partial \beta} = 2X^T X \beta - 2X^T y + 2\lambda \beta$$
 Set the equation to 0, and we get: $2X^T X \beta + 2\lambda \beta = 2X^T y$ $\beta = (X^T X + \lambda I)^{-1} X^T y$

1.2

No if we don't apply regularization, because $X^T X = \begin{bmatrix} 1 & 3 & 5 \\ 2 & 6 & 10 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 3 & 6 \\ 5 & 10 \end{bmatrix} = \begin{bmatrix} 35 & 70 \\ 70 & 140 \end{bmatrix}$, which is not invertible.

If we use Ridge regression, however, the question becomes solvable, because $\begin{bmatrix} 35 & 70 \\ 70 & 140 \end{bmatrix} - \lambda I$ will be invertible

1.3

Lasso regression, because it prefers 0 coefficients

2.2 Programming

```
1 # Machine Learning HW2 Ridge Regression
2
3 import matplotlib.pyplot as plt
4 import numpy as np
5
6 # Parse the file and return 2 numpy arrays
7
8
9 def load_data_set(filename):
10     arr = np.loadtxt(filename)
11     x, y = np.split(arr, [-1], axis=1)
12     return x, y
13
14 # Split the data into train and test examples by the train_proportion
15 # i.e. if train_proportion = 0.8 then 80% of the examples are training and 20%
16 # are testing
17
18
19 def train_test_split(x, y, train_proportion):
20     num_train = int(x.shape[0] * train_proportion)
21     x_train, x_test = x[:num_train, :], x[num_train:, :]
22     y_train, y_test = y[:num_train, :], y[num_train:, :]
23     return x_train, x_test, y_train, y_test
24
25 # Find theta using the modified normal equation
26 # Note: lambdaV is used instead of lambda because lambda is a reserved word in python
27
28
29 def normal_equation(x, y, lambdaV):
30     # your code
31     x_t = np.transpose(x)
32     n = x.shape[1]
33     beta = np.dot(np.dot(np.linalg.inv(
34         np.dot(x_t, x) + lambdaV * np.identity(n)), x_t), y)
35     return beta
36
37 # Extra Credit: Find theta using gradient descent
38
39
40 def gradient_descent(x, y, lambdaV, num_iterations, learning_rate):
41     # your code
42     return beta
43
44 # Given an array of y and y_predict return loss
45
46
47 def get_loss(y, y_predict):
48     diff = y - y_predict
49     loss = np.dot(diff.T, diff) / len(y)
50     return loss[0][0]
51
52 # Given an array of x and theta predict y
53
54
55 def predict(x, theta):
56     # your code
57     y_predict = x.dot(theta)
```

```

58     return y_predict
59
60 # Find the best lambda given x_train and y_train using 4 fold cv
61
62
63 def cross_validation(x_train, y_train, lambdas, n_folds=4):
64     valid_losses = []
65     training_losses = []
66     num_examples = x_train.shape[0]
67     num_per_fold = num_examples//n_folds
68     # your code
69     for lambda_ in lambdas:
70         valid_loss_sum = 0
71         training_loss_sum = 0
72         for i in range(n_folds):
73             testing_start = i*num_per_fold
74             beta = normal_equation(np.concatenate((x_train[0:testing_start],
x_train[testing_start+num_per_fold:]), axis=0), np.concatenate(
75                 (y_train[0:testing_start], y_train[testing_start+num_per_fold:]), axis=0),
lambda_)
76             training_loss_sum += get_loss(np.concatenate((y_train[0:testing_start],
y_train[testing_start+num_per_fold:]), axis=0), predict(
77                 np.concatenate((x_train[0:testing_start], x_train[testing_start+num_per_fold:]),
axis=0), beta))
78             valid_loss_sum += get_loss(y_train[testing_start: testing_start+num_per_fold],
predict(
79                 x_train[testing_start: testing_start+num_per_fold], beta))
80             valid_losses.append(valid_loss_sum/n_folds)
81             training_losses.append(training_loss_sum/n_folds)
82
83     return np.array(valid_losses), np.array(training_losses)
84
85
86 def bar_plot(best_beta):
87     x = range(1, best_beta.shape[0]+1)
88     plt.bar(x=x, height=best_beta.flatten())
89     plt.title("Final beta bar graph")
90     plt.show()

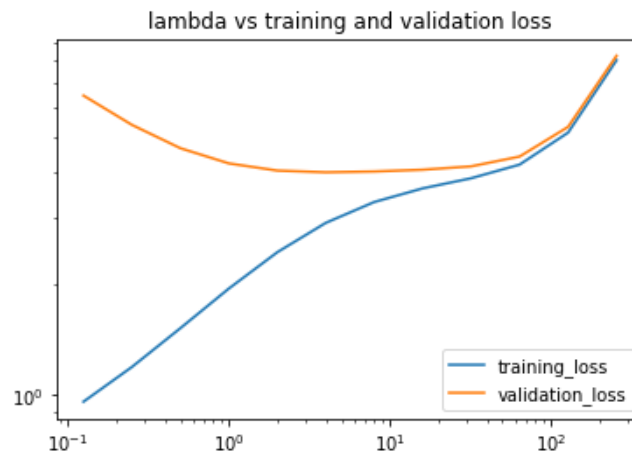
```

Task1

```

1  # step 1
2  # If we don't have enough data we will use cross validation to tune hyperparameter
3  # instead of a training set and a validation set.
4  x, y = load_data_set("dataRidge.txt") # load data
5  x_train, x_test, y_train, y_test = train_test_split(x, y, 0.8)
6  # Create a list of lambdas to try when hyperparameter tuning
7  lambdas = [2**i for i in range(-3, 9)]
8  lambdas.insert(0, 0)
9  # Cross validate
10 valid_losses, training_losses = cross_validation(x_train, y_train, lambdas)
11 # Plot training vs validation loss
12 plt.plot(lambdas[1:], training_losses[1:], label="training_loss")
13 # exclude the first point because it messes with the x scale
14 plt.plot(lambdas[1:], valid_losses[1:], label="validation_loss")
15 plt.legend(loc='best')
16 plt.xscale("log")
17 plt.yscale("log")
18 plt.title("lambda vs training and validation loss")
19 plt.show()
20
21 best_lambda = lambdas[np.argmin(valid_losses)]
22 print('best lambda: {}'.format(best_lambda))

```



```
1 best_lambda: 4
```

As shown in the above graph, training loss keeps increasing as lambda increases because it prevents the regression from overfitting the training data. Validation loss, on the other hand, decreases when we increase lambda, that is, up to lambda=4, because it not overfitting the test data set allows the model to generalize better. But if we increase the value of lambda beyond 4, the model starts to have a hard time learning the pattern in the data, so the validation loss goes up.

Task2

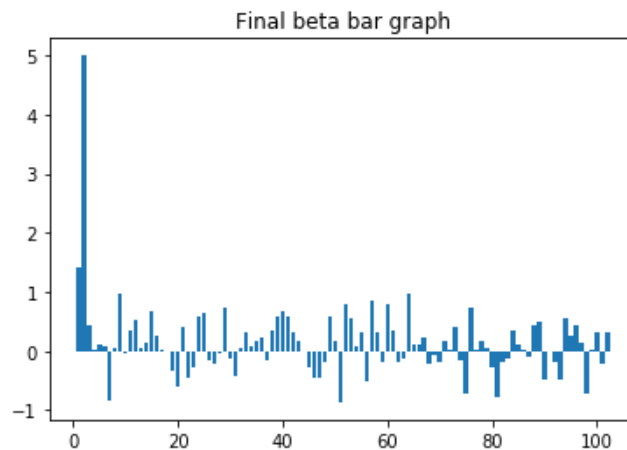
```
1 # step 2: analysis
2 normal_beta = normal_equation(x_train, y_train, 0)
3 best_beta = normal_equation(x_train, y_train, best_lambda)
4 large_lambda_beta = normal_equation(x_train, y_train, 512)
5 # your code get l2 norm of normal_beta
6 normal_beta_norm = np.linalg.norm(normal_beta)
7 # your code get l2 norm of best_beta
8 best_beta_norm = np.linalg.norm(best_beta)
9 # your code get l2 norm of large_lambda_beta
10 large_lambda_norm = np.linalg.norm(large_lambda_beta)
11 print('best lambda: {}'.format(best_lambda))
12 print("L2 norm of normal beta: " + str(normal_beta_norm))
13 print("L2 norm of best beta: " + str(best_beta_norm))
14 print("L2 norm of large lambda beta: " + str(large_lambda_norm))
15 print("Average testing loss for normal beta: " +
16       str(get_loss(y_test, predict(x_test, normal_beta))))
17 print("Average testing loss for best beta: " +
18       str(get_loss(y_test, predict(x_test, best_beta))))
19 print("Average testing loss for large lambda beta: " +
20       str(get_loss(y_test, predict(x_test, large_lambda_beta))))
21
22 # Step3: Retrain a new model using all sampling in training, then report error on testing set
23 # your code !
24 final_beta = normal_equation(x_train, y_train, best_lambda)
25 print("Final testing loss: " +
26       str(get_loss(y_test, predict(x_test, final_beta))))
27
28 # Step Extra Credit: Implement gradient descent, analyze and show it gives the same or very
29 # similar beta to normal_equation
30 # to prove that it works
```

```
1 best_lambda: 4
2 L2 norm of normal beta: 30.269654798800733
3 L2 norm of best beta: 6.640092215033666
4 L2 norm of large lambda beta: 4.651284674384335
5 Average testing loss for normal beta: 11.031286619643865
6 Average testing loss for best beta: 4.636831026037389
7 Average testing loss for large lambda beta: 12.126420332837494
8 Final testing loss: 4.636831026037389
```

As lambda increases, we get smaller norms of β , which helps generalization, but up to a point. If the lambda is too large, the model won't be able to properly fit the data

Task3

```
1 bar_plot(final_beta)
```



The data generation function is $5x + 3 + \text{noise}$, so overall the model was able to capture the pattern of the data fairly well.

Sample questions

Question 1. Basis functions for regression

No. This is not a good selection of basis functions in that all the three basis function don't overlap in the range the x where they have effect, i.e. ϕ_1 only has effect when $0 \leq x \leq 2$, ϕ_2 only has effect when $2 \leq x \leq 4$, ϕ_3 only has effect when $4 \leq x \leq 6$, rendering the basis functions unable to complement one another, and giving the model little flexibility. For example, this model can't fit a straight line

Question 2. Polynomial Regression

0, because $y = x^2$ fits the data perfectly If we leave (1, 1) out, we get $4 = 2\beta_1 + 4\beta_2$ and $9 = 3\beta_1 + 9\beta_2$, which gives us $\beta_2 = 1, \beta_1 = 0$ This gives us an error of $1^2 - 1 = 0$

If we leave (2, 4) out, we get $1 = \beta_1 + \beta_2$ and $9 = 3\beta_1 + 9\beta_2$, which also gives us $\beta_2 = 1, \beta_1 = 0$ This gives us an error of $2^2 - 4 = 0$

Leaving (3, 9) out, we get $4 = 2\beta_1 + 4\beta_2$ and $1 = \beta_1 + \beta_2$, which, again, gives us $\beta_2 = 1, \beta_1 = 0$ This gives us an error of $3^2 - 9 = 0$ So the mse = $(0 + 0 + 0)/3 = 0$