# ECE469 Homework 2

Chase A. Lotito

October 5, 2024

## Question 1

**Solution.**

q1.py

```python
###########################
# Chase Lotito - SIUC F24 #
# ECE469 - Intro to ML    #
# HW2 - Question 1        #
###########################

# IMPORT LIBRARIES
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.preprocessing import PolynomialFeatures
from sklearn.model_selection import cross_val_score,
    train_test_split, KFold
from sklearn.linear_model import LinearRegression
from sklearn.linear_model import Ridge
from sklearn.metrics import mean_squared_error

# Read in provided csv data to pandas dataframe
RAW_DATA_PATH = 'C:/Users/cloti/OneDrive/Desktop/CODE/datasets/
    datasetHW2P1.csv'
data = pd.read_csv(RAW_DATA_PATH)

test_size = 0.2


# (A) SPLIT DATASET TO CREATE TWO SUB DATASETS FOR TRAINING AND
    TESTING
x = data['x'].values
y = data['y'].values

x = x.reshape(-1, 1)
```

```python
29
30 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=
      test_size, random_state=42)
31
32
33
34
35 # (B) USE POLYNOMIAL REGRESSION TO FIT 5 POLYNOMIAL MODELS (DEG. 1 -
      DEG. 5)
36 deg = [1, 2, 3, 4, 5]
37 mse_train_arr = []
38 mse_test_arr = []
39
40 for i in deg:
41     # transform input features into polynomial features
42     poly = PolynomialFeatures(degree=i)    # initialize polynomial
43     x_train_poly = poly.fit_transform(x_train)
44     x_test_poly = poly.transform(x_test)
45
46     # fit model
47     poly_regressor = LinearRegression()
48     poly_regressor.fit(x_train_poly, y_train)
49
50     # test
51     y_train_predicted = poly_regressor.predict(x_train_poly)
52     y_test_predicted = poly_regressor.predict(x_test_poly)
53
54     # calc mse
55     mse_train = mean_squared_error(y_train, y_train_predicted)
56     mse_test = mean_squared_error(y_test, y_test_predicted)
57
58     # add mse's to arrays for later plotting
59     mse_train_arr.append(mse_train)
60     mse_test_arr.append(mse_test)
61
62     # print mse for training and testing
63     print(f"Degree: {i}")
64     print(f"MSE Training: {mse_train}")
65     print(f"MSE Testing: {mse_test}")
66
67 # plotting mse for training and testing against model complexity
68 plt.plot(deg, mse_train_arr, label='Training')
69 plt.plot(deg, mse_test_arr, label='Testing')
70 plt.xlabel('Model Complexity (Degree)')
71 plt.ylabel('Mean Square Error')
72 plt.title('Model Error v. Model Complexity')
73 plt.legend()
74 #plt.show()
```

```python
75  #plt.close()
76
77
78
79
80  # (C) USE 10-FOLD CROSS-VALIDATION TO FIND THE MODEL WHICH OPTIMALLY
         FITS GIVEN DATASET.
81  #       PLOT THE TRAINING, CROSS-VALIDATION, AND TESTING ERRORS
         AGAINST MODEL COMPLEXITY.
82
83  # parameters
84  k = 10                      # for 10-fold cross-validation
85  best_degree = 1         # track best degree
86  max_degree = 5
87  best_mse = float('inf')
88  mse_per_degree = []
89
90  # perform k-fold cross-validation for each polynomial degree
91  for degree in range(1, max_degree + 1):
92      kf = KFold(n_splits=k, shuffle=True, random_state=42)
93      mse_fold_values = []
94      k_iter = 1
95
96      # initialize a plot for regression fit
97      all_y_pred = np.zeros(x.shape)
98
99      # Loop through each fold
100     for train_index, test_index in kf.split(x):
101         x_train, x_test = x[train_index], x[test_index]
102         y_train, y_test = y[train_index], y[test_index]
103
104         # transform original x data into polynomial features for
                 current degree
105         poly = PolynomialFeatures(degree=degree)    # here degree is
                 the parent loop iterator
106         x_train_poly = poly.fit_transform(x_train)
107         x_test_poly = poly.transform(x_test)
108
109         # intialize and fit linear regression on polynomial space
                 input features
110         poly_regressor = LinearRegression()
111         poly_regressor.fit(x_train_poly, y_train)         # <--
                 TRAINING MODEL HERE
112
113         # predict using trained model
114         y_test_pred = poly_regressor.predict(x_test_poly)
115
116         # calc mse for test set
```

```python
117             mse_test = mean_squared_error(y_test, y_test_pred)
118             mse_fold_values.append(mse_test)
119
120         # calculuate avg mse for all folds for current deg
121         avg_mse = np.mean(mse_fold_values)
122         mse_per_degree.append(avg_mse)
123
124         # check if current deg has lowest avg mse
125         if avg_mse < best_mse:
126             best_mse = avg_mse
127             best_degree = degree
128
129
130         # Generate new plot for degree
131         plt.figure(figsize=(12,8))
132         plt.scatter(x,y,s=10,color='blue', label='Original Data', alpha
                =0.6)
133
134         # plot polynomial fit for degree
135         sorted_x = np.sort(x, axis=0)
136         sorted_x_poly = poly.transform(sorted_x)
137         plt.plot(sorted_x, poly_regressor.predict(sorted_x_poly), color=
                'red', label=f"Degree {degree}", alpha = 0.7)
138
139         # finalize plot for degree
140         plt.title(f"Polynomial Regression (Degree {degree})")
141         plt.xlabel('x')
142         plt.ylabel('y')
143         plt.legend(loc='best')
144         plt.grid(True)
145         plt.savefig(f"poly_reg_deg{degree}_fold{k_iter}.png", dpi=300)
146         plt.close()
147
148 # print best degree and mse
149 print(f"Best Polynomial Degree: {best_degree}")
150
151 # plot the mse v. polynomial degree
152 plt.figure(figsize=(10,6))
153 plt.plot([1,2,3,4,5], mse_per_degree, marker='o', color='b', label='
        Cross-Validation MSE', alpha=0.75)
154 plt.plot(deg, mse_train_arr, label='Training MSE', marker='+', alpha
        =0.75)
155 plt.plot(deg, mse_test_arr, label='Testing MSE', marker='x', alpha
        =0.75)
156 plt.title("MSE vs Model Complexity")
157 plt.xlabel('Polynomial Degree')
158 plt.ylabel('Mean Squared Error (MSE)')
159 plt.xticks(range(1, max_degree + 1))
```

```python
160  plt.legend()
161  plt.grid(True)
162  plt.savefig('plots\\mse_vs_complexity.png', dpi=300)
163  plt.close()
164
165
166
167
168  # (D) CONSIDER A 4-DEGREE POLYNOMIAL AS YOUR MODEL.
169  #      USE RIDGE REGRESSION AND FIND BEST HYPERPARAMETER
170  #      LAMBDA VIA 10-FOLD CROSS-VALIDATION. PLOT THE CROSS
171  #      VALIDATION ERROR VERSUS LN(LAMBDA).
172
173  # remember we have x as input features and y as output features
174
175  # parameters
176  degree = 4                            # polynomial degree
177  lambdas = np.logspace(-4, 2, 100)   # lambdas for ridge
178  kf = KFold(n_splits=10, shuffle=True, random_state=1)
179
180  # map input features into its polynomial space
181  poly = PolynomialFeatures(degree=degree)
182  x_poly = poly.fit_transform(x)
183
184  # store cross-validation results
185  mse_values = []
186
187  # perform cross-validation for each lambda
188  for alpha in lambdas:
189      ridge = Ridge(alpha=alpha)
190
191      # calc mse using cross_val_score
192      mse = -cross_val_score(ridge, x_poly, y, cv=kf, scoring='
             neg_mean_squared_error').mean()
193      mse_values.append(mse)
194
195  # find lambda with best minimum mse
196  best_lambda = lambdas[np.argmin(mse_values)]
197  print(f"Best lambda (w/ minimum MSE): {best_lambda}")
198
199  # fit model with best lambda
200  ridge_best = Ridge(alpha=best_lambda)
201  ridge_best.fit(x_poly, y)
202
203  # generate testing x values
204  x_fit = np.linspace(x.min(), x.max(), 1000).reshape(-1,1)
205  x_poly_fit = poly.transform(x_fit)    # send x_fit to polynomial
         space
```

```
206  y_fit = ridge_best.predict(x_poly_fit)
207
208  # plot mse vs. ln(lambda)
209  plt.figure(figsize=(8,6))
210  plt.plot(np.log(lambdas), mse_values, label='MSE')
211  plt.xlabel('$\\log_e (\\lambda)$')
212  plt.ylabel('Cross-Validated MSE')
213  plt.title('Cross-Validated MSE v. $\\log_e (\\lambda)$')
214  plt.grid(True)
215  plt.legend()
216  plt.savefig('crossvalmse_vs_loglambda.png', dpi=300)
217  plt.close()
```

```
Degree: 1
MSE Training: 19458.420066787396
MSE Testing: 19657.90065574125
Degree: 2
MSE Training: 12843.093582390631
MSE Testing: 12953.786640428427
Degree: 3
MSE Training: 5.011946346822575
MSE Testing: 5.000928251849794
Degree: 4
MSE Training: 5.011930151180137
MSE Testing: 5.000749381529411
Degree: 5
MSE Training: 5.011819378630573
MSE Testing: 5.000656140106037
Best Polynomial Degree: 3
Best lambda (w/ minimum MSE): 0.24770763559917114
```
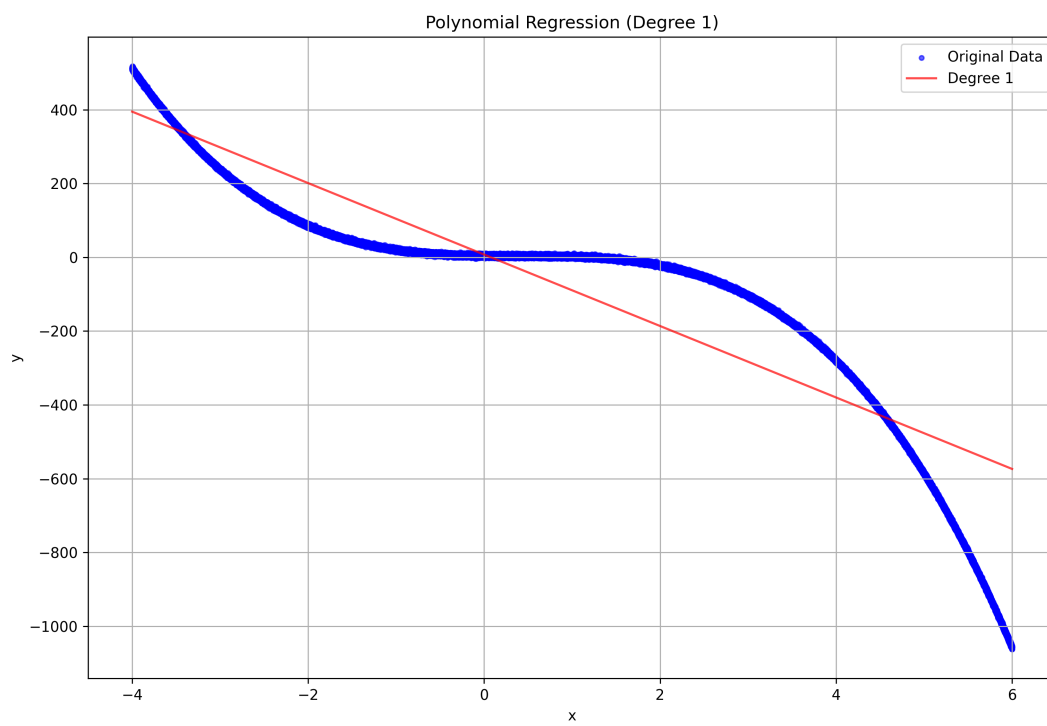
Figure 1: q1.py Terminal Output

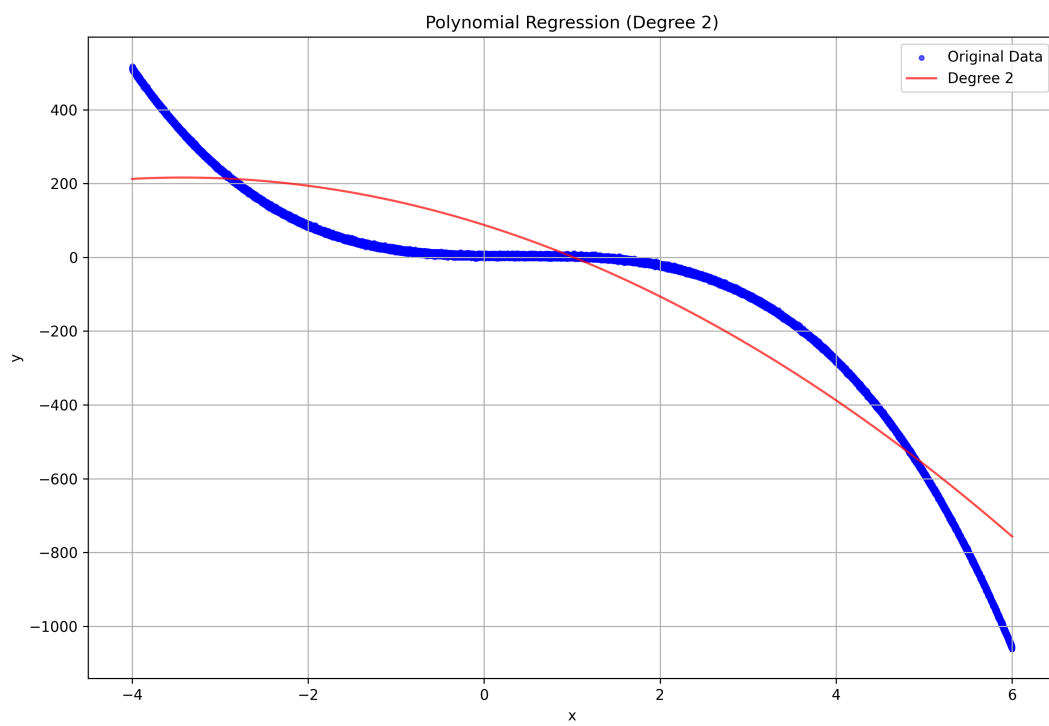Figure 2: Polynomial Model Degree 1
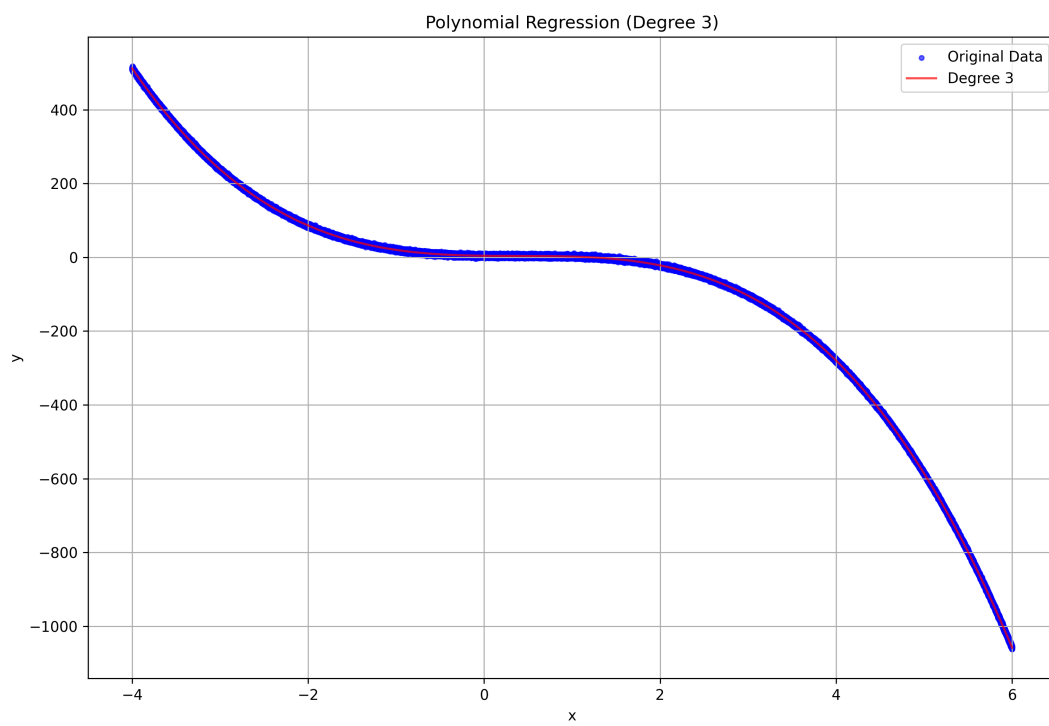
Figure 3: Polynomial Model Degree 2
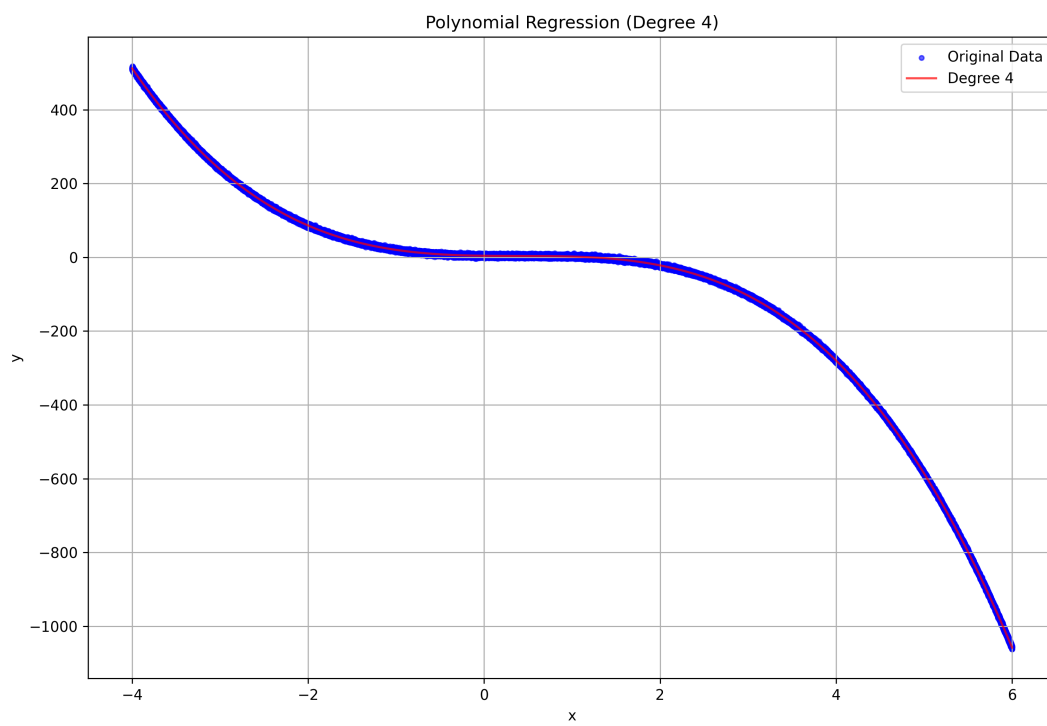
Figure 4: Polynomial Model Degree 3

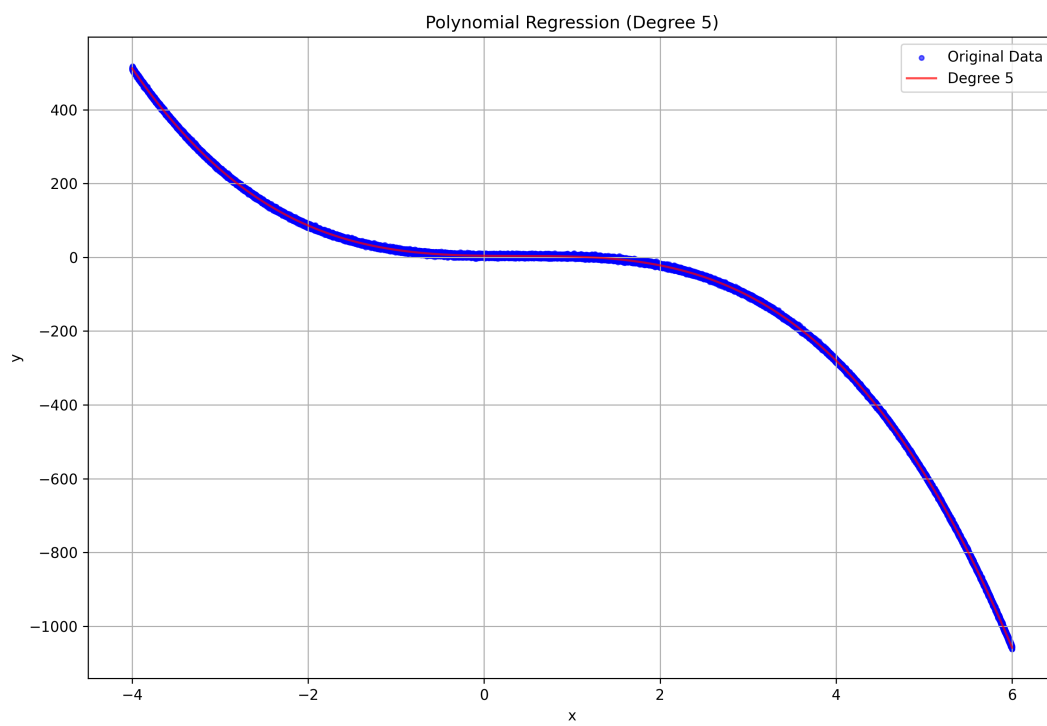Figure 5: Polynomial Model Degree 4

Figure 6: Polynomial Model Degree 5
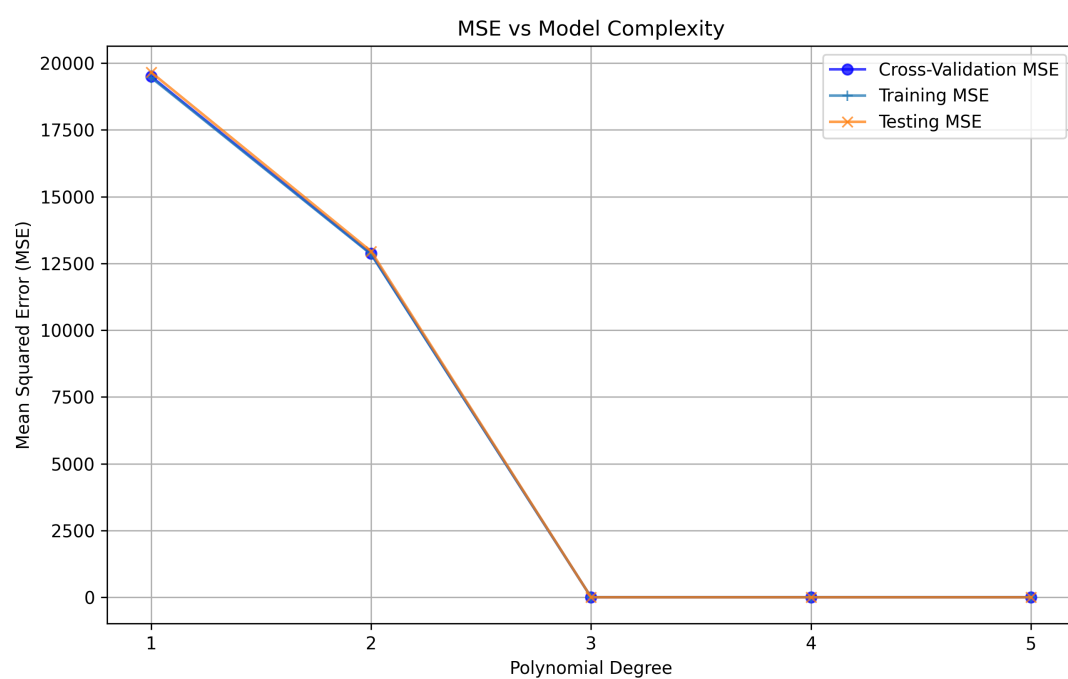
Figure 7: MSE v. Model Complexity

Figure 8: Cross-Validated MSE v. $\log_e \lambda$

## Question 2

**Solution.**

q2.py

```python
###########################
# Chase Lotito - SIUC F24 #
# ECE469 - Intro to ML    #
# HW2 - Question 2        #
###########################

# IMPORT LIBRARIES
import numpy as np
import pandas as pd
from sklearn.preprocessing import OrdinalEncoder    # For encoding
    categorical features
from sklearn.impute import SimpleImputer            # For adding
    missing values
from sklearn.preprocessing import StandardScaler    # For
    standardizing data
from sklearn.preprocessing import PolynomialFeatures
from sklearn.model_selection import cross_val_score,
    train_test_split, KFold
from sklearn.linear_model import LinearRegression
from sklearn.linear_model import Ridge
from sklearn.metrics import mean_squared_error
import matplotlib.pyplot as plt


# (A) DOWNLOAD HOUSING.CSV

# Get housing data
RAW_DATA = 'https://github.com/ageron/data/raw/main/housing/housing.
    csv'
housing = pd.read_csv(RAW_DATA)



# (B) DATA-PREPROCESSING FROM HW1

# Choose input features and output features (saved into numpy.
    ndarray type)
X = housing[
        ['longitude',
        'latitude',
        'housing_median_age',
        'total_rooms',
```

```python
37              'total_bedrooms',
38              'population',
39              'households',
40              'median_income',
41              'ocean_proximity']
42        ].values
43   Y = housing[['median_house_value']].values
44
45   # Ocean Proximity is a categorical feature. Drop it or transform
         into numerical values (encode).
46
47   # Isolate the ocean_proximity data in input data X
48   ocean_proximity = X[:,8].reshape(-1,1)    # reshape(-1,1) to make 2D
         array for Ordinal
49   # Initalize the ordinal encoder
50   ordinal_encoder = OrdinalEncoder()
51   # Encode the ocean_proximity strings into numerical data
52   encoded_ocean = ordinal_encoder.fit_transform(ocean_proximity)
53   # Put the encoded version of ocean_proximity into input data X
54   X[:,8] = encoded_ocean.flatten()           # flatten to add 1D version
         of array back into X
55
56   # Clean the dasta by either dropping or replacing missing values
57
58   # Initialized SimpleImputer, will use the median to add missing
         entries
59   simple_imputer = SimpleImputer(strategy='median')
60
61   # Change X np ndarray into a Pandas Dataframe to use SimpleImputer
62   dX = pd.DataFrame(X)
63   dY = pd.DataFrame(Y)
64
65   # Perform SimpleImputer transformation, for both inputs X and
         outputs Y
66   imputed_data = simple_imputer.fit_transform(dX)
67   X = imputed_data
68   imputed_data = simple_imputer.fit_transform(dY)
69   Y = imputed_data
70
71   # Carry out feature scaling either via normalization or
         standardization.
72   std_scaler = StandardScaler()
73   scaled_data = std_scaler.fit_transform(X)
74   X = scaled_data
75   scaled_data = std_scaler.fit_transform(Y)
76   Y = scaled_data
77
78   # set test set size
```

```python
79  test_size = 0.2
80
81  # split into testing and training set (both outputted as pd.
        DataFrames)
82  X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=
        test_size, random_state=42)
83
84
85
86
87  # (C) USE LINEAR REGRESSION TO DEVELOP AN ML MODEL FOR PREDICTION OF
88  #      'MEDIAN_HOUSE_VALUE' FOR FUTURE INPUTS AND ANALYZE TEST ERRORS
89  #        EXPLICITLY EXPRESS THE CORRESPONDING OPTIMAL WEIGHTS AND THE
90  #        FINAL LEARNED MODEL. USE GRAPHICAL REPRESENTATIONS.
91
92  # initalize and train linear model
93  linear_regressor = LinearRegression()
94  linear_regressor.fit(X_train, Y_train)   # <-- train model here
95
96  # extract the optimal weights from the ML model
97  weights = linear_regressor.coef_.flatten()      # flatten makes it a
        normal list
98
99  # test model
100 Y_train_predicted = linear_regressor.predict(X_train)
101 Y_test_predicted = linear_regressor.predict(X_test)
102
103 # calculate mean square error
104 mse_train = mean_squared_error(Y_train, Y_train_predicted)
105 mse_test = mean_squared_error(Y_test, Y_test_predicted)
106
107 # VISUALIZATION
108 # print out results of linear model
109 print("--------------------")
110 print("LINEAR MODEL RESULTS")
111 print("--------------------")
112 print(f"Model Weights: {weights}")
113 print(f"MSE (train): {mse_train*100:.2f}%")
114 print(f"MSE (test): {mse_test*100:.2f}%")
115
116 # Plot predicted vs actual values
117 plt.scatter(Y_test, Y_test_predicted, marker='o', s=0.75, c='#32a852
        ', alpha=0.95, label='Pred v. Actual')     # plot predicted
        against actual, if diagonalized, well fit.
118 plt.plot([min(Y_test), max(Y_test)], [min(Y_test), max(Y_test)],
        label='Ideal Model',color='red', linewidth=2)
119 plt.xlabel("Actual Median House Value")
120 plt.ylabel("Predicted Median House Value")
```

```python
121 plt.title("Predicted Median House Value vs Actual Median House Value
        ")
122 plt.legend()
123 plt.savefig("plots\\q2_pred_vs_actual.png", dpi=300)
124 plt.close()
125
126 # Plot the Learned Weights
127 # Coefficients of the model
128 #feature_names = ['longitude', 'latitude', 'housing_median_age', '
        total_rooms',
129 #                 'total_bedrooms', 'population', 'households', '
        median_income', 'ocean_proximity']
130 #
131 #plt.barh(feature_names, weights)
132 #plt.xlabel("Model Weights $w_i$")
133 #plt.title("Linear Regression Feature Weights")
134 #plt.show()
135 #plt.close()
136
137
138
139 # (D) USE CROSS-VALIDATION TECHNIQUES TO IMPROVE THE GENERALIZATION
        OF THE MODEL AND ANALYZE
140 #     THE ROOT MEAN SQUARE ERROR (RMSE). USE GRAPHICAL ILLUSTRATIONS
        .
141
142 # parameters
143 lambdas = np.logspace(-4, 2, 100)  # lambdas for ridge
144 kf = KFold(n_splits=10, shuffle=True, random_state=1)
145
146 # store cross-validation results
147 mse_values = []
148
149 # perform cross-validation for each lambda
150 for alpha in lambdas:
151     ridge = Ridge(alpha=alpha)
152
153     # calc mse using cross_val_score
154     mse = -cross_val_score(ridge, X, Y, cv=kf, scoring='
            neg_mean_squared_error').mean()
155     mse_values.append(mse)
156
157 # find lambda with best minimum mse
158 best_lambda = lambdas[np.argmin(mse_values)]
159 print("--------------------")
160 print("RIDGE REGULARIZATION")
161 print("--------------------")
162 print(f"Best lambda (w/ minimum MSE): {best_lambda:.2f}")
```

```
163  print(f"Minimum MSE: {np.min(mse_values)*100:.2f}%")
164
165  # fit model with best lambda
166  ridge_best = Ridge(alpha=best_lambda)
167  ridge_best.fit(X_train, Y_train)
168
169  # testing x values
170  Y_test_ridge_pred = ridge_best.predict(X_test)
171
172  # plot rmse vs. ln(lambda)
173  plt.figure(figsize=(8,6))
174  plt.plot(np.log(lambdas), np.sqrt(mse_values), c="#570710",label='
         MSE')
175  plt.xlabel('$\\log_e (\\lambda)$')
176  plt.ylabel('Cross-Validated RMSE')
177  plt.title('Q2: Cross-Validated RMSE v. $\\log_e (\\lambda)$')
178  plt.grid(True)
179  plt.legend()
180  plt.savefig('plots\\q2_crossvalmse_vs_loglambda.png', dpi=300)
181  plt.close()
```

We can plot, like in Figure 10, the predicted values against the actual values to see how well the learned model works. If the predicted values (green) lie along the same line as the actual values (red), then the model is ideal. However, we see the linear model has a large spread, which means the model is most likely not complex enough to model the input features.



```
--------------------
LINEAR MODEL RESULTS
--------------------
Model Weights: [-0.72960647 -0.77646856  0.12237936 -0.14310869  0.3233905  -0.3668318
  0.22907347  0.65294944  0.00340332]
MSE (train): 36.25%
MSE (test): 37.30%
--------------------
RIDGE REGULARIZATION
--------------------
Best lambda (w/ minimum MSE): 14.17
Minimum MSE: 36.55%
```

Figure 9: q2.py Terminal Output

From Figure 9, we can write the equation for the final learned model:

$$y(\mathbf{x}) = -0.7296x_1 + -0.7765x_2 + 0.1224x_3 + -0.1431x_4 + 0.3234x_5$$
$$+ -0.3668x_6 + 0.2291x_7 + 0.6529x_8 + 0.0034x_9$$

Where the input vector $\mathbf{x}$ is defined as:

18

$$\mathbf{x} = \begin{bmatrix} x_1: & \text{longitude} \\ x_2: & \text{latitude} \\ x_3: & \text{housing median age} \\ x_4: & \text{total rooms} \\ x_5: & \text{total bedrooms} \\ x_6: & \text{population} \\ x_7: & \text{households} \\ x_8: & \text{median income} \\ x_9: & \text{ocean proximity} \end{bmatrix}$$
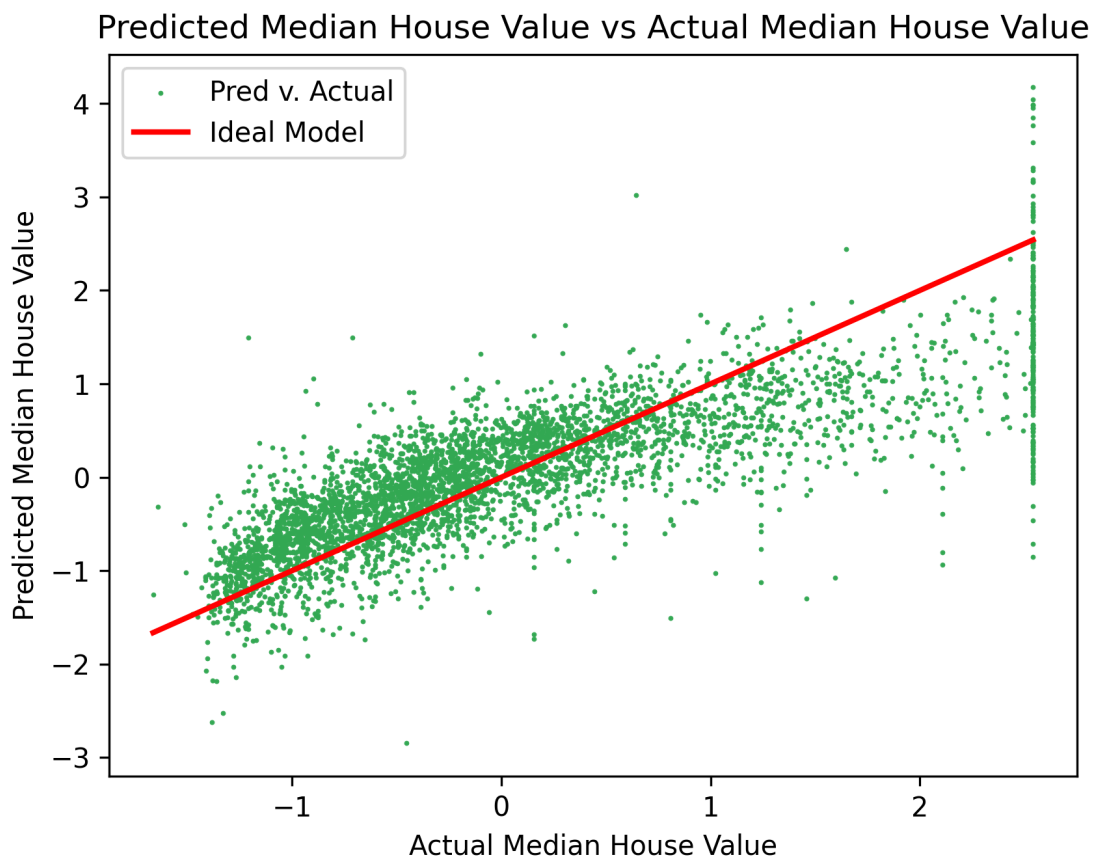


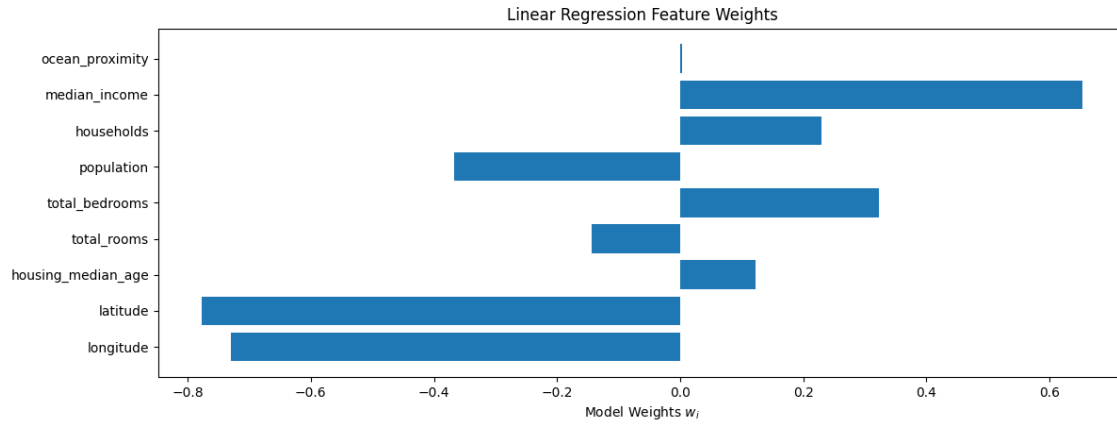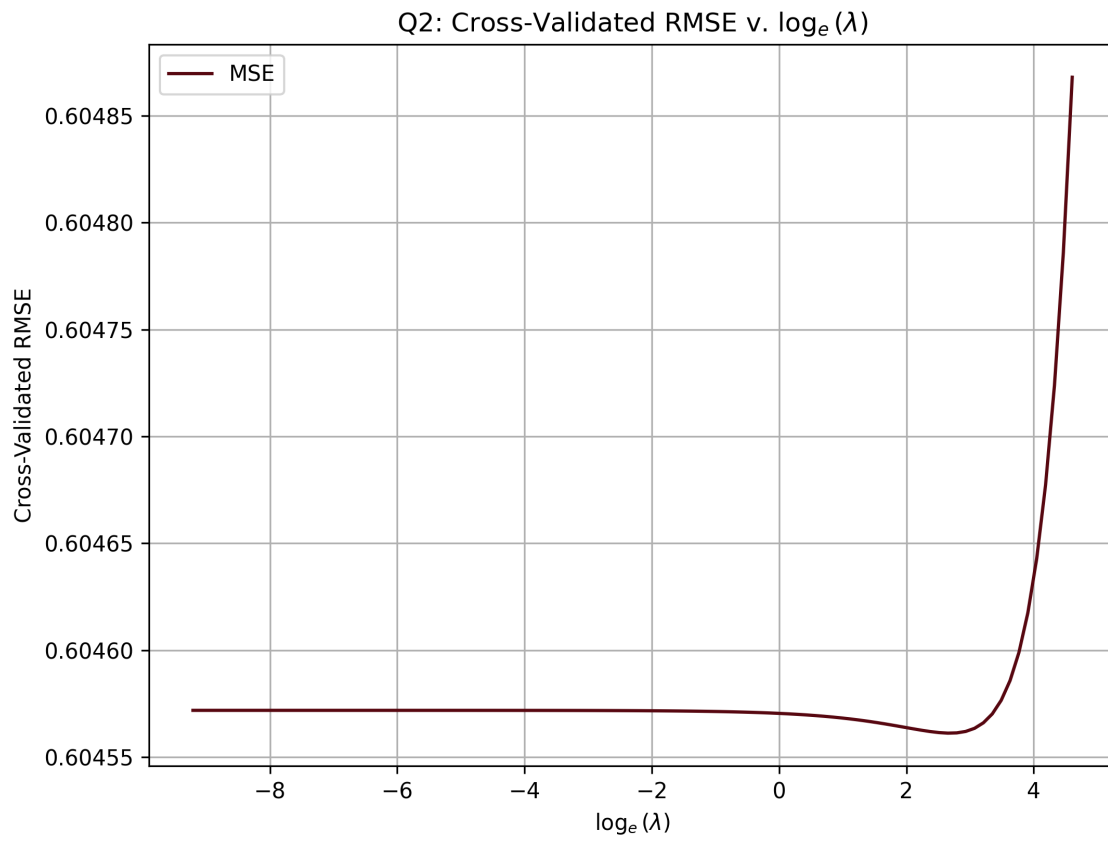Figure 10: Predicted Median House Value v. Actual Median House Value

Figure 11: Model Weights



Figure 12: Cross-Validated RMSE v. $\log_e \lambda$