

Supervised Learning (COMP0078) Report

Yoga Advaith Veturi zchayav@ucl.ac.uk

Maki Egawa zcabmeg@ucl.ac.uk

November 16, 2020

1 PART I

The source code for this section is provided in `p1_sourcecode.ipynb`, which is attached with the submission, or is accessible via this clickable [link](#). For convenience however, relevant code snippets are also attached inline with the explanations.

1.1 Linear Regression

- (a) We are provided with a toy data set

$$\mathcal{D} = \{(1, 3), (2, 2), (3, 0), (4, 5)\}$$

where each pair (x_i, y_i) contains an input $x_i \in \mathbb{R}^n$ and output $y_i \in \mathbb{R}$. The goal is to fit \mathcal{D} using a linear combination of polynomial basis functions

$$\{\phi_1(x) = 1, \phi_2(x) = x^2, \phi_3(x) = x^3, \dots, \phi_k(x) = x^{k-1}\}$$

of dimension k , where each basis function $\phi_k : \mathbb{R}^n \rightarrow \mathbb{R}$. The fitting is performed for dimensions $k = 1, 2, 3, 4$.

To perform the curve fitting, we first defined a function Φ that takes the vector of inputs \mathbf{x} and transforms it into an $(m \times k)$ feature map $\Phi(\mathbf{x})$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} \longrightarrow \Phi(\mathbf{x}) = \begin{bmatrix} \phi_1(x_1) & \dots & \phi_k(x_1) \\ \phi_1(x_2) & \dots & \phi_k(x_2) \\ \vdots & \ddots & \vdots \\ \phi_1(x_m) & \dots & \phi_k(x_m) \end{bmatrix}$$

where each $\phi_k(x)$ is a single polynomial basis function. The following code executes this, provided the value of k .

```

1 #Computes the polynomial basis
2 def phi(x, k):
3     """ Computes a single basis function  $x^{k-1}$ , given
4         an input vector x of shape (m, 1) and the value k.
5         m is the number of training examples """
6     return x**(k-1)
7
8 #This function creates the feature map PHI()
9 def PHI(x, k):
10     """ Creates the feature map, which is a matrix of
11         shape (m, k). Each column is the kth basis function
12         applied on
13         the input vector x of shape (m, 1). m is the number
14         of training examples """
15
16     #Number of training examples
17     m = x.shape[0]
18
19     #Initialize a feature map of shape (m, k) with 0s in
20     it.
21     feature_map = np.zeros((m, k))
22
23     #In every column, replace the zeros with the basis
24     function k computed on input vector x
25     for i in range(0, k):
26         feature_map[:, i] = phi(x.squeeze(), i+1)
27
28     return feature_map
29
30 #Compute feature maps for k = {1,...,4}
31 x_tr_k1 = PHI(x, k=1)

```

```

25 x_tr_k2 = PHI(x, k=2)
26 x_tr_k3 = PHI(x, k=3)
27 x_tr_k4 = PHI(x, k=4)

```

With the transformed feature map $\Phi(\mathbf{x})$, we then computed the weights $\mathbf{w} \in \mathbb{R}^k$ that minimises the sum of squared errors defined by

$$SSE = \sum_{i=1}^m (y_i - \mathbf{w} \cdot \phi(x))^2 = (\Phi \mathbf{w} - y)^T (\Phi \mathbf{w} - y)$$

Setting the gradient $\nabla_{\mathbf{w}}(SSE) = 0$ and solving for \mathbf{w} , the least squares solution is

$$\mathbf{w} = (\Phi^T \Phi)^{-1} \Phi^T y$$

The following code computes the weights for $k = 1, 2, 3, 4$.

```

1 #Compute weights for the transformed data using least
  squares regression
2 w_k1 = np.linalg.inv(x_tr_k1.T @ x_tr_k1) @ x_tr_k1.T @
  y
3 w_k2 = np.linalg.inv(x_tr_k2.T @ x_tr_k2) @ x_tr_k2.T @
  y
4 w_k3 = np.linalg.inv(x_tr_k3.T @ x_tr_k3) @ x_tr_k3.T @
  y
5 w_k4 = np.linalg.inv(x_tr_k4.T @ x_tr_k4) @ x_tr_k4.T @
  y

```

- (b) The following polynomial equations were obtained, corresponding to $k = 1, 2, 3, 4$:

$$y = 2.5$$

$$y = 1.5 + 0.4x$$

$$y = 9.0 - 7.1x + 1.5x^2$$

$$y = -5.0 + 15.17x - 8.5x^2 + 1.33x^3$$

These are visualized in Figure 1.

- (c) For each fitted curve, the mean squared error was computed using $MSE = \frac{SSE}{m}$. The following code executes the MSE function.

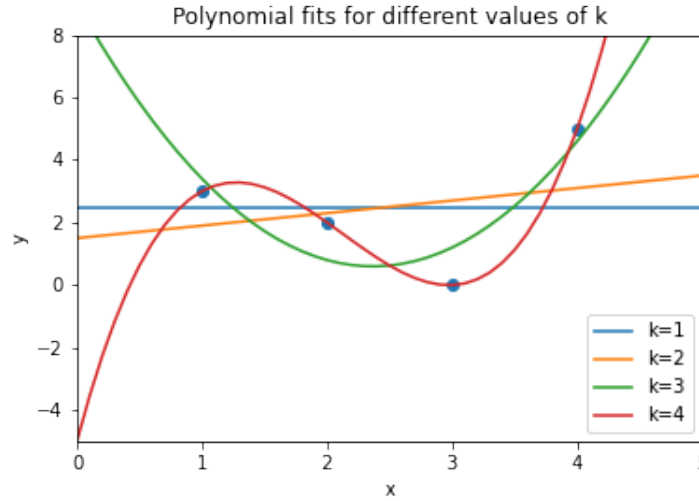


Figure 1: Fitted curves for $k=1,2,3,4$

```

1 def MSE(x_values, y_values, weights):
2     """
3     Computes the MSE given the x_values matrix (m, n),
4     y_values vector of shape (m, 1) and weights vector
5     of shape (n, 1).
6     m is the number of training examples, n is the number
7     of input features.
8     """
9     m = x_values.shape[0]
10    SSE = ((x_values @ weights) - y_values).T @ ((
11        x_values @ weights) - y_values)
12    MSE = SSE / m
13    return MSE
14
15 #Print MSE values for each of the polynomial bases
16 print("MSE for k = 1 : " + str(MSE(x_tr_k1, y, w_k1).
17     item()))
18 print("MSE for k = 2 : " + str(MSE(x_tr_k2, y, w_k2).
19     item()))
20 print("MSE for k = 3 : " + str(MSE(x_tr_k3, y, w_k3).
21     item()))
22 print("MSE for k = 4 : " + str(MSE(x_tr_k4, y, w_k4).
23     item()))

```

The MSE values are presented in Table 1 below.

k	MSE
1	3.25
2	3.05
3	0.80
4	3.49e-23

Table 1: Mean Squared Error for fitted curves $k = 1, 2, 3, 4$

2. We now fit models on data generated from the random function

$$g_{\sigma=0.07}(x) := \sin^2 2\pi x + \epsilon$$

where ϵ is Gaussian noise with a mean of 0 and variance $\sigma^2 = 0.07^2$.

- (a) i. A training set of 30 points was generated by sampling random points uniformly in the interval $[0, 1]$. These points were passed through g_σ and on each call of the function, Gaussian noise ϵ is added. The following code implements this:

```

1 def g_sigma(x, sigma):
2     """
3     Random function of form g_sigma(x) = sin^2(2*pi*x
4     ) + epsilon where epsilon is a Gaussian
5     distributed random variable
6     with mean 0 and sigma^2 variance.
7
8     Takes an input x and returns the value g_sigma(x)
9     with an added noise value
10    generated from epsilon.
11    """
12
13    #Compute function on the point x
14    function = np.square(np.sin(2 * np.pi * x))
15    #Generate a single noise value for that point x
16    noise = np.random.normal(0.0, sigma)
17
18    return function + noise
19
20 def generate_data(interval, size, sigma):
21     """
22     Generates a random sample {(x_1, g_1), ..., (x_m,
23     g_m)}. First samples a single x value from the
24     interval

```

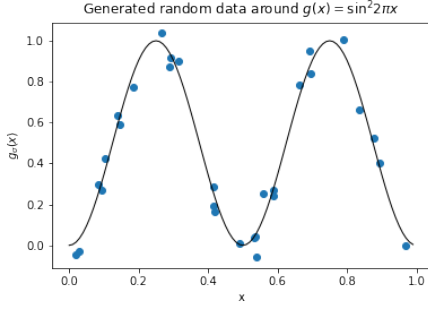
```

20 and then computes  $g\_sigma = \sin^2(2\pi x) +$ 
   epsilon, where epsilon is random gaussian noise.
21
22 Args
23 ----
24 interval - tuple of (lower_bound, upper_bound)
25 size - number of points to generate
26 sigma - standard deviation of the noise
27
28 Returns
29 -----
30 NumPy arrays of the generated input and output
   data
31 """
32
33 #This will store the x_values and g_sigma(x)
   values
34 x_values = []
35 g_values = []
36
37 for i in range(size):
38
39     #Generate a single point
40     x = np.random.uniform(interval[0], interval[1])
41     #Create the output value using the random
   function
42     g = g_sigma(x, sigma=sigma)
43     #Append the values to the list
44     x_values.append(x)
45     g_values.append(g)
46
47     return np.array(x_values), np.array(g_values)
48
49 #Generate training data - set a seed for ensuring
   reproducible results
50 np.random.seed(1)
51 x_train, g_train = generate_data(interval=(0, 1),
   size=30, sigma=0.07)

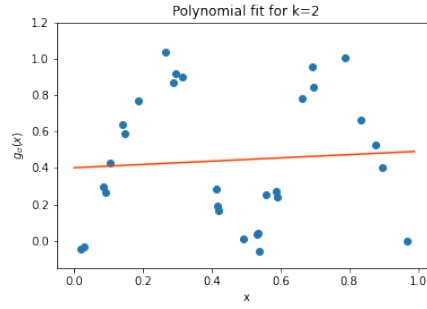
```

The generated data set is plotted along with the deterministic function $\sin^2(2\pi x)$ in Figure 2a.

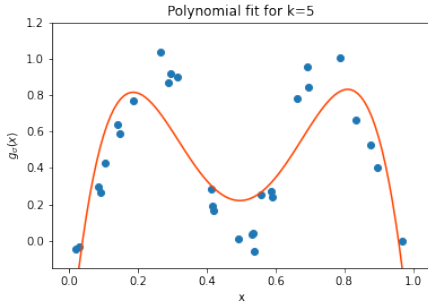
- ii. This data set is then fit using polynomial bases of dimension $k = 2, 5, 10, 14, 18$ using the same process followed in 1a. The x values are first transformed into the feature map Φ , then using



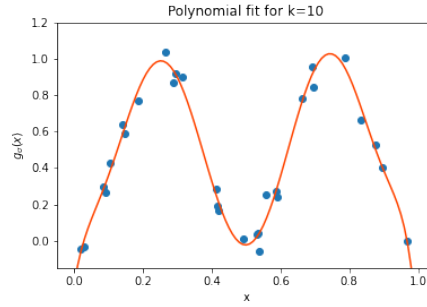
(a)



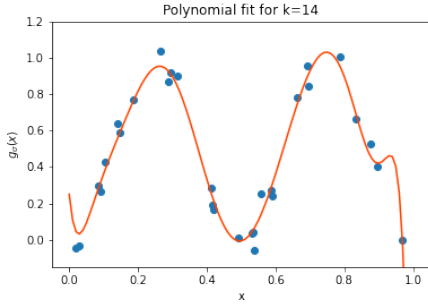
(b)



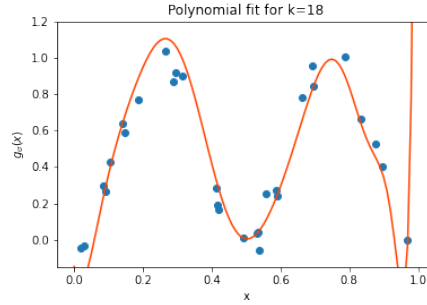
(c)



(d)



(e)



(f)

Figure 2: (a) Generated data (blue points) with the deterministic function $\sin^2(2\pi x)$ overlapped (black curve). All fitted curves are shown in orange. (b) Fitted curve for $k = 2$. (c) Fitted curve for $k = 5$. (d) Fitted curve for $k = 10$. (e) Fitted curve for $k = 14$. (f) Fitted curve for $k = 18$.

least squares regression, the weights \mathbf{w} are computed for each k . These fitted curves are presented in Figure 2(b-f).

- (b) The training error denoted te_k was computed by performing linear regression on the training points for every value of $k = 1, \dots, 18$ and calculating the MSE on the training points using the $MSE()$ function. The natural log of the training error is plotted in Figure 3(a). It is noticeable that the function is decreasing, however around epoch 16, there is a rise, which could be attributed to numerical errors.
- (c) A test set of 1000 points was generated using the *generate_data()* function described earlier. The test set error, denoted tse_k was computed by calculating the MSE on the test points using the weights computed on the training points. The natural log of the test set error for $k = 1, \dots, 18$ is plotted in Figure 3(b).
- (d) The previous results correspond to one type of generated training and test set. To see these errors smoothed out for different types of training and test sets, the previous steps in 2b and 2c were repeated for 100 runs and the log of the average train and test set error over the runs were computed. The results are plotted in Figure 3(c).

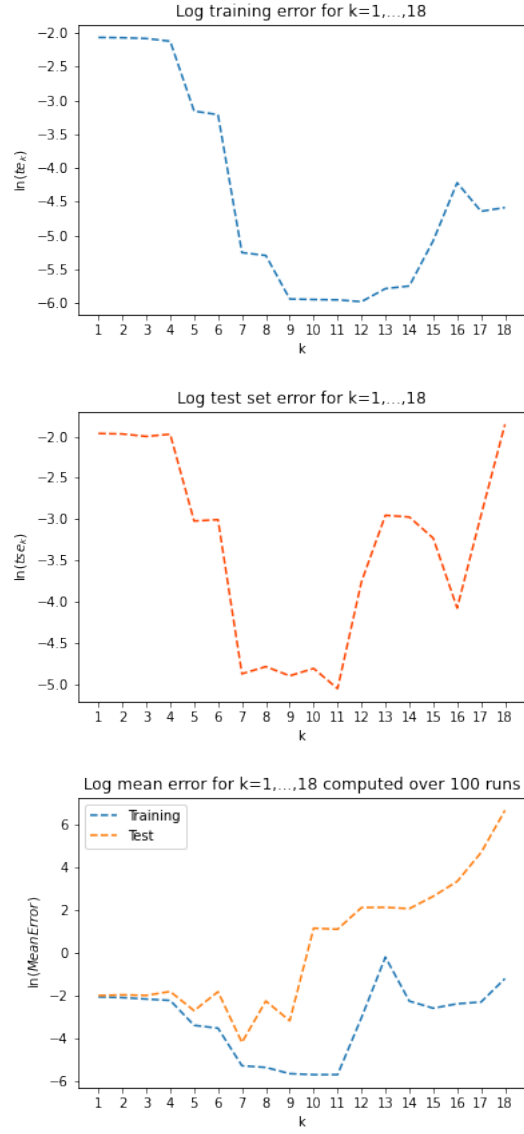


Figure 3: MSE values for $k = 1, \dots, 18$ (a) Log training error (b) Log test set error (c) Log training and test set error, averaged over 100 runs

3. The experiments performed in 2b, 2c and 2d are repeated, now using a different basis

$$\{\sin 1\pi x, \sin 2\pi x, \dots, \sin k\pi x\}$$

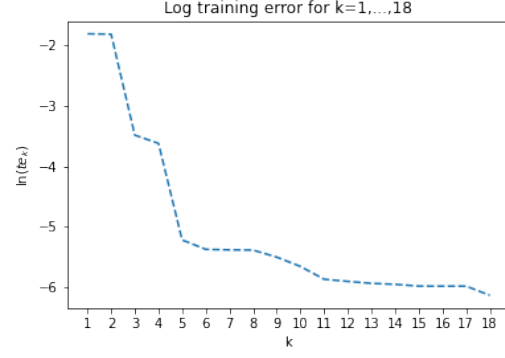
We modify the functions for computing the feature map as follows:

```

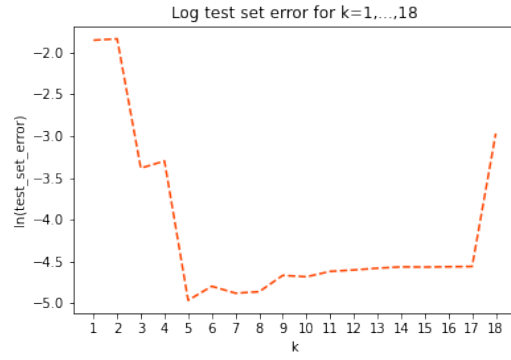
1 #Computes the basis function
2 def phi_v2(x, k):
3     """ Computes a single basis function, given input vector
4       x of shape (m, 1) and value k. m is the number of
5       training examples """
6     return np.sin(k*np.pi*x)
7
8 #This function creates the feature map
9 def PHI_v2(x, k):
10    """ Creates the feature map, which is a matrix of shape (
11      m, k). Each column is the kth basis function applied on
12      the input vector x of shape (m, 1). m is the number of
13      training examples """
14
15    #Extract the number of training examples
16    m = x.shape[0]
17
18    #Initialize feature map matrix to matrix of zeros
19    feature_map = np.zeros((m, k))
20
21    #Iterative over each column
22    for i in range(0, k):
23        #For each column of the feature map matrix, replace the
24        #0s with the computed kth basis function values
25        feature_map[:, i] = phi_v2(x.squeeze(), i+1)
26    return feature_map

```

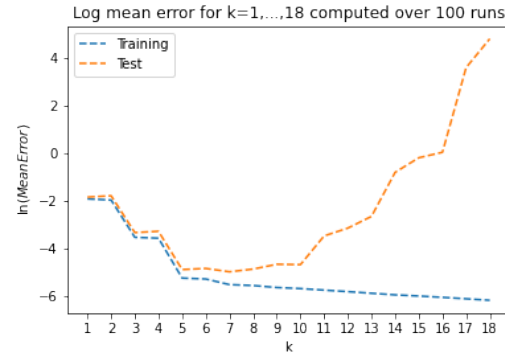
The log te_k and tse_k for $k = 1, \dots, 18$ are plotted in Figure 4a and 4b respectively. The entire operation is performed for 100 runs, using different generated training and test points. The log of the average te_k and tse_k is visualized in Figure 4c.



(a)



(b)



(c)

Figure 4: (a) Log training error for $k = 1, \dots, 18$ (b) Log test error for $k = 1, \dots, 18$ (c) Log of average train and test error over 100 runs.

1.2 Filtered Boston Housing and Kernels

4. In this section, we were provided with a modified version of the Boston Housing data set. Below is a quick inspection of the first few rows of this data.

CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	LSTAT	MEDV
0.00632	18.0	2.31	0	0.538	6.575	65.2	4.0900	1	296	15.3	4.98	24.0
0.02731	0.0	7.07	0	0.469	6.421	78.9	4.9671	2	242	17.8	9.14	21.6
0.02729	0.0	7.07	0	0.469	7.185	61.1	4.9671	2	242	17.8	4.03	34.7
0.03237	0.0	2.18	0	0.458	6.998	45.8	6.0622	3	222	18.7	2.94	33.4
0.06905	0.0	2.18	0	0.458	7.147	54.2	6.0622	3	222	18.7	5.33	36.2

Table 2: Inspection of first 5 rows of Boston Housing data set

The goal was to predict the median house price (MEDV) using the remaining 12 variables. To do so, we performed multiple experiments with different variants of linear regression. These experiments involved splitting the data into a training set, which contained 2/3 of the examples, and a test set with 1/3 of the examples. The following function implements the data splitting. Note that the data can be shuffled before splitting in order to generate random splits.

```
1 def split_data(inputs, targets, test_proportion, shuffle=
    None):
2     """
3     Splits the data into training and test sets.
4
5     Args
6     ----
7     inputs : NumPy array of input data. Should be of shape (#
        examples, # features)
8     targets : NumPy array of target data. Should be of shape
        (# examples, 1)
9     test_proportion : Value between 0 and 1 which specifies
        how much of the data to use for testing.
10    shuffle : Optional. Set to True if you want the data
        shuffled and then split.
11    seed : Optional. Set for reproducible results.
12
13    Returns
14    -----
15    train_X : NumPy array of training examples. Should be of
        shape (# examples, # features)
```

```

16 train_Y : NumPy array of training targets. Should be of
    shape (# examples, 1)
17 test_X : NumPy array of testing examples. Should be of
    shape (# examples, # features)
18 test_Y : NumPy array of testing targets. Should be of
    shape (# examples, 1)
19 """
20
21 #Stores the number of data points
22 nData = inputs.shape[0]
23
24 # Shuffle data
25 if shuffle:
26     #Generate a shuffled version of the array indices
27     shuffled_indices = np.random.permutation(nData)
28     #Shuffle the inputs as per in the array of shuffled
    indices
29     shuffled_inputs = inputs[shuffled_indices, :]
30     shuffled_targets = targets[shuffled_indices, :]
31 else:
32     #If shuffle is set to False then we just work with the
    data in its original order
33     shuffled_inputs = inputs
34     shuffled_targets = targets
35
36 # Calculate the split index based on the specified
    proportions
37 split_index = int((1 - test_proportion) * nData)
38
39 # Select the examples up to the split index to be used as
    training set
40 train_X = shuffled_inputs[:split_index]
41 train_Y = shuffled_targets[:split_index]
42 # Select the examples from the split index onwards to be
    used as the test set
43 test_X = shuffled_inputs[split_index:]
44 test_Y = shuffled_targets[split_index:]
45
46 return train_X, train_Y, test_X, test_Y

```

- (a) We first performed Naive regression. The data set was split into train and test sets using the splitting function. A vector of ones was created with the same length as the training set and another for the test set. Using the training ones vector and the training outputs, we computed the weights with linear regression. Finally, the train and test MSE was computed using the $MSE()$ function. The code implementation for this is provided below. Note that it is performed over 20 runs for 20 random train/test splits.

```
1 #For reproducibility of results
2 np.random.seed(101010)
3
4 #Stores the number of runs
5 runs = 20
6
7 #These will store the total train and test MSE over the
   20 runs
8 train_MSE = 0
9 test_MSE = 0
10
11 #Iterate over runs
12 for i in range(runs):
13
14     #Generate samples
15     xtrain, ytrain, xtest, ytest = split_data(inputs=X,
        targets=Y, test_proportion=1/3, shuffle=True)
16
17     train_ones = np.ones((xtrain.shape[0], 1))
18     test_ones = np.ones((xtest.shape[0], 1))
19
20     #Perform Naive Regression
21     w_naive = np.linalg.inv(train_ones.T @ train_ones) @
        train_ones.T @ ytrain
22
23     #Calculate and aggregate train and test MSE
24     train_MSE += MSE(train_ones, ytrain, w_naive)
25     test_MSE += MSE(test_ones, ytest, w_naive)
26
27     print(">> Run {} Completed.".format(i+1))
```

The average train and test MSE over 20 runs is as follows:

```
Average Train MSE = 86.05624289198637
Average Test MSE = 81.51650453132758
```

- (b) With Naive Regression, we are essentially fitting the data using a constant function and asking the question "What would the model predict in the absence of any features?". The interpretation of this constant function is as follows:

We can define the training inputs and outputs as the vectors X_{train} and Y_{train}

$$X_{train} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}, Y_{train} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix}$$

where m is the number of training examples. If we used these vectors to compute the weights w_{naive} then we would get the following solution:

$$w_{naive} = (X_{train}^T X_{train})^{-1} X_{train}^T Y_{train} = \frac{1}{m} \sum_{i=1}^m y_i$$

Hence, the naive weights is essentially the mean of the training output values, which is what the model would predict for a given input. To validate this, we compared the output of the w_{naive} with the mean of Y_{train} on a single run and found them to be the same:

```
w_naive value = [[23.11780415]]
Mean of Y_train values = 23.11780415430267
```

- (c) We next performed Linear Regression over single attributes. A column of ones was added to the matrix of data in order to account for the bias term in the linear regression. The following function adds the ones column:

```

1 #Function to add ones to the training set as a way to
  account for the bias term in the weights
2 def add_bias(data):
3     """
4     Adds a column of ones into the dataset in order to
      account for the bias term.
5
6     Args
7     ----
8     data : a NumPy array of the input data. Can be a
      matrix of shape (m, n) or a single input vector x_i
      of shape (m, 1). m is the number of training
      examples
9     while n is the number of input features.
10
11    Returns
12    -----
13    data_with_bias : a NumPy array of data with 1s added.
      If data is a matrix, then a column of 1s is added
      to the first column.
14    If data is a single vector x_i, then the first
      element of x_i is 1.
15    """
16
17    #If the data is a vector
18    if data.shape[1] == 1:
19        #Create a one
20        ones = np.ones((data.shape[1], ))
21        #Insert 1 as the first value for the vector
22        data_with_bias = np.insert(data, 0, ones, axis=0)
23
24    #If data is a matrix
25    else:
26        #Create a vector ones
27        ones = np.ones((data.shape[0], ))
28        #Stack this vector of ones to the beginning of the
      matrix
29        data_with_bias = np.insert(data, 0, ones, axis=1)
30
31    return data_with_bias
32
33 #Convert the pandas dataframe into matrix
34 data = dataframe.values
35
36 #Add a column of 1's to the data matrix to account for

```



```

    bias term
37 data_with_bias = add_bias(data)
38
39 #Split the data matrix into a matrix of the inputs and
    the vector of outputs
40 X = data_with_bias[:, :13]
41 Y = data_with_bias[:, 13].reshape(-1, 1)

```

Using the *split_data()* function, the data was shuffled and then split into training and testing sets. For training, the individual columns of the training data matrix, which represent the different attributes, were selected along with the "ones" column and the weights were computed using linear regression for these "sub-matrices". The training and test MSE was then computed using the *MSE()* function. We report the results averaged over 20 runs for 20 random train/test splits in Table 3 below.

Method	MSE Train	MSE Test
Attribute 1	69.568609	77.953024
Attribute 2	72.015690	76.825691
Attribute 3	62.422693	69.677222
Attribute 4	80.442598	85.280471
Attribute 5	67.203539	73.004103
Attribute 6	42.845189	45.544419
Attribute 7	70.498032	76.720579
Attribute 8	77.077643	83.831789
Attribute 9	70.034179	76.850144
Attribute 10	63.711957	70.728458
Attribute 11	61.489686	65.364829
Attribute 12	37.487598	40.817604

Table 3: Train and Test MSE for Linear regression over single attributes. Averaged over 20 runs.

- (d) The next experiment was linear regression over all attributes together. We followed the exact same procedure as in the single attributes, except now we computed the weights using the entire training data matrix. We report the training and test MSE averaged over

20 runs for 20 random train/test splits as follows.

```
Average Train MSE = 21.414764727427052
Average Test MSE  = 25.922665880996654
```

It can be seen that our train and test MSE for linear regression over all attributes is smaller than the regressors for single attributes. This is expected as our models incorporated more features, allowing for better prediction of the output values.

1.3 Kernelized Ridge Regression

5. In this section, we perform kernel ridge regression over the training set with ℓ examples. This involved computing the dual solution $\boldsymbol{\alpha}^*$ that optimises the equation

$$\boldsymbol{\alpha}^* = \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^\ell} \frac{1}{\ell} \sum_{i=1}^{\ell} \left(\sum_{j=1}^{\ell} \alpha_j K_{i,j} - y_i \right)^2 + \gamma \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha} \quad (1)$$

where $K(x_i, x_j)$ (abbreviated $K_{i,j}$) is the Gaussian kernel computed using x_i and x_j :

$$K(x_i, x_j) = \exp \left(-\frac{\|x_i - x_j\|^2}{2\sigma^2} \right) \quad (2)$$

\mathbf{K} is an $(\ell \times \ell)$ matrix containing the entries $K_{i,j}$.

The dual solution is thus defined as

$$\boldsymbol{\alpha}^* = (\mathbf{K} + \gamma \ell \mathbf{I}_\ell)^{-1} \mathbf{y} \quad (3)$$

where γ is a regularization parameter, \mathbf{I}_ℓ is an $(\ell \times \ell)$ identity matrix and \mathbf{y} is the vector of training outputs.

To evaluate the regression function on a test point, the formula is

$$y_{test} = \sum_{i=1}^{\ell} \alpha_i^* K(x_i, x_{test}) \quad (4)$$

- (a) In order to perform kernel ridge regression, the data set was split into training and testing sets (2/3:1/3) using the splitting function. As can be seen in equations 1 and 2, there are two additional hyperparameters σ and γ . In order to select the best values for these hyperparameters, we first performed kernel ridge regression on the training set using five-fold cross-validation. The chosen range of values for σ and γ were $[2^7, 2^{7.5}, \dots, 2^{13}]$ and $[2^{-40}, 2^{-39}, \dots, 2^{-26}]$ respectively.

The following are the steps for 5-fold cross-validation on a single pair of (σ, γ) . We performed this over all pairs of hyperparameters.

- i. Shuffle the training set. (To reduce bias)
- ii. Divide into $K=5$ groups.
- iii. Select first group as the "validation" set and use the remaining $K-1$ groups as the training set.
- iv. Compute the dual solution $\boldsymbol{\alpha}^*$ using the training set.
- v. Compute predictions for validation data.
- vi. Compute MSE for validation set and record value.
- vii. Repeat steps (iii-vi) on all groups.

In order to perform steps (iv) and (v), three components were crucial: the computation of the kernel matrix \mathbf{K} , the dual solution $\boldsymbol{\alpha}^*$ and the evaluation function (4). The code for these three components are presented below. Note that with the computation of \mathbf{K} and the evaluation function, the simple method would be to use for-loops. However, this can be computationally expensive, hence we developed a vectorized implementation, for which a detailed explanation is provided in Appendix 3.1.

```

1 def gaussian_kernel_matrix(X, sigma):
2     """
3     Computes the gaussian kernel matrix on a given set of
4     vectors.
5
6     Args
7     ----
8     X : NumPy array of data. Should be of shape (l, n)
9         where l is number of examples and n is number of
10        features
11    sigma : sigma value of Gaussian kernel.
12
13    Returns
14    -----
15    K : a symmetric positive definite matrix of shape (l,
16        l)
17    """
18
19    #Number of samples
20    l = X.shape[0]
21
22    # Extract the diagonals of  $XX^T$  into a vector and
23    # stack it l times into lxl matrix B
24    inner_product_vector = np.diagonal(X@X.T).reshape(-1,
25        1)
26
27    #Compute matrix of just the dot products between x_i
28    # and x_j
29    inner_product_matrix = inner_product_vector - 2 * (
30        X@X.T) + inner_product_vector.T
31
32    #Compute gaussian kernel matrix
33    K = np.exp((-1/(2*sigma**2)) * inner_product_matrix)
34
35    #Checkpoint to ensure matrix K is correct dimensions
36    assert K.shape == (l, l), "K matrix should be of
37        shape ({}, {})".format(l, l)
38
39    return K
40
41 def kernel_ridge_regression(X, y, sigma, gamma):
42     """
43     Compute the solution  $\alpha^*$  for kernelized ridge
44     regression.
45

```

```

36  Args
37  ----
38  X : NumPy array of data. Should be of shape (l, n)
      where l is number of examples and n is number of
      input features
39  y : NumPy array of target outputs. Should be of shape
      (l, 1)
40  sigma : sigma value for Gaussian kernel
41  gamma : regularization parameter for ridge regression
42
43  Returns
44  -----
45  alpha_star : (l, 1) vector which represents the dual
      solution for ridge regression.
46
47  """
48
49  #Extract number of samples
50  l = X.shape[0]
51
52  #Compute kernel
53  K = gaussian_kernel_matrix(X, sigma) #kernel_matrix(X
      , X, sigma=sigma)
54
55  #Create identity matrix used in the formula
56  I = np.eye(l)
57
58  #Compute alpha_star value using formula
59  try:
60      alpha_star = np.linalg.inv(K + gamma*l*I) @ y
61  except np.linalg.LinAlgError:
62      alpha_star = np.linalg.pinv(K + gamma*l*I) @ y
63
64  #Checkpoint to ensure the alpha is correct dimensions
65  assert alpha_star.shape == (l, 1), "Alpha star should
      be of shape ({}, {})".format(l, 1)
66
67  return alpha_star
68
69  def evaluate(alpha, X_tr, X_te, sigma):
70      """
71      Vectorized implementation of kernel ridge regression
      evaluation function on a test point.
72
73      Args

```

```

74  ----
75  alpha : NumPy array of the dual solution alpha.
       Should be of shape (l, 1) where l is the number of
       training examples.
76  X_tr : NumPy array of the training set inputs. Should
       be of shape (l, n)
77  X_te : NumPy array of the test set inputs. Should be
       of shape (m, n) where m is number of test examples.
78  sigma : Sigma value of the Gaussian kernel.
79
80  Returns
81  -----
82  pred : NumPy array of predictions for the test data.
       Should be of shape (m, 1)
83  """
84
85  #Extract dimensions of data
86  l = X_tr.shape[0]
87  m = X_te.shape[0]
88
89  #Evaluate regression model at all test points
90  inner_product_matrix = np.diagonal(X_tr@X_tr.T).
       reshape(1, -1) - 2 * (X_te @ X_tr.T) + np.diagonal(
       X_te@X_te.T).reshape(-1, 1)
91  K = np.exp((-1/(2*sigma**2)) * inner_product_matrix)
92  pred = K @ alpha
93
94  assert pred.shape == (m, 1), "Predictions should be
       of shape ({}, 1)".format(m, 1)
95
96  return pred

```

The final implementation of the entire 5-fold cross-validation is as follows:

```

1  def train_with_kfoldCV(k, data, kernel_params, shuffle=
       True, verbose=True):
2      """
3      Performs k-fold cross validation on a given dataset.
4
5      Args
6      ----
7      k - number of folds of cross-validation
8      data - tuple of training input and output data.
       Training input must be NumPy array of shape (#
       examples, # features) while output must be of shape

```

```

9      (# examples, 1)
10     kernel_params - tuple of (sigma, gamma) parameters
11                     for computing the Gaussian kernel.
12     shuffle - If True, then data is shuffled
13     seed - to ensure reproducible results
14     verbose - Set to True if you wish to see printed
15               notifications.
16
17     Returns
18     -----
19     cv_error - the validation set error averaged over the
20                 k folds of cross validation for the specified
21                 kernel parameters
22     """
23
24     #Extract data
25     x, y = data
26
27     #Extract gaussian kernel params
28     sigma, gamma = kernel_params
29
30     #Extract dimensions
31     nSamples, nFeatures = x.shape
32
33     #Shuffle dataset randomly
34     if shuffle:
35         perm = np.random.permutation(nSamples)
36         x_shuffled = x[perm, :]
37         y_shuffled = y[perm, :]
38     else:
39         x_shuffled = x
40         y_shuffled = y
41
42     # print("X_shuffled = ", x_shuffled)
43
44     #Split data into k-groups
45     x_groups = np.array_split(x_shuffled, k)
46     y_groups = np.array_split(y_shuffled, k)
47
48     #Stores the cross validation MSE which is the MSE
49         calculated over the validation sets generated during
50         each fold of k-fold CV
51     cv_mse = 0
52
53     #Iterate over each cross-validation group

```



```

47 for i in range(len(x_groups)):
48
49     #Use the selected group as "validation" set
50     val_inputs, val_outputs = x_groups[i], y_groups[i]
51
52     #Use rest of groups as training set
53     train_inputs = np.vstack([x_groups[j] for j in
54 range(len(x_groups)) if j != i])
54     train_outputs = np.vstack([y_groups[j] for j in
55 range(len(x_groups)) if j != i])
56
57     #Fit model on train set
58     alpha_star = kernel_ridge_regression(X=train_inputs
59 , y=train_outputs, sigma=sigma, gamma=gamma)
60
61     #Generate predictions on "validation" set
62     val_predictions = evaluate(alpha=alpha_star, X_tr=
63 train_inputs, X_te=val_inputs, sigma=sigma)
64
65     #Compute validation errors
66     cv_mse += (1/len(val_outputs)) * (val_outputs -
67 val_predictions).T @ (val_outputs - val_predictions)
68     if verbose:
69         print(">> Cross-validation Fold {} Completed.".
70 format(i+1))
71
72 #Average the errors
73 cv_error = cv_mse.item() / k
74
75 if verbose:
76     print(">> MEAN CROSS-VALIDATION ERROR = {}\n".
77 format(cv_error))
78
79 #Summarize performance using sample of model eval
80 scores
81 return cv_error
82
83 #Create parameter space
84 sigma_values = 2*np.arange(7, 13.5, 0.5)
85 gamma_values = (2*np.arange(40.0, 25.0, -1))**-1
86
87 #Represent all (sigma, gamma) pairs using meshgrid.
88 This will be helpful when creating the 3D plot
89 gamma_array, sigma_array = np.meshgrid(gamma_values,
90 sigma_values)

```

```

82
83 #This grid will store the cross validation error values
84 errors = np.zeros((sigma_values.shape[0], gamma_values.
85                    shape[0]))
86
87 #Split data into train and test set - for
88 reproducibility, I've set a seed
89 np.random.seed(100000)
90
91 X_train, Y_train, X_test, Y_test = split_data(inputs=X
92        [:, 1:], targets=Y, test_proportion=1/3, shuffle=
93         True)
94
95 train_data = (X_train, Y_train)
96
97 #Train the model for all the different values of sigma
98 and gamma
99
100 for i in range(errors.shape[0]):
101     for j in range(errors.shape[1]):
102
103         #Create the sigma and gamma pair
104         sigma = sigma_values[i]
105         gamma = gamma_values[j]
106
107         print("Sigma = {}, Gamma = {}".format(sigma, gamma)
108               )
109
110         #Perform 5-fold cross validation. Note - Because we
111         are adding a shuffle operation inside crossval
112         function, every pair of sigma and gamma may train on
113         different examples
114         #To ensure that each (sigma, gamma) pair is used on
115         the same set of training and test examples on each
116         fold, I've seed a seed
117         np.random.seed(130399)
118         cv_error = train_with_kfoldCV(k=5, data=train_data,
119                                     kernel_params=(sigma, gamma), shuffle=True, verbose
120                                     =True)
121
122         #Add this error value into our errors matrix
123         errors[i, j] = cv_error
124
125 print("Ridge Regression Completed!")

```

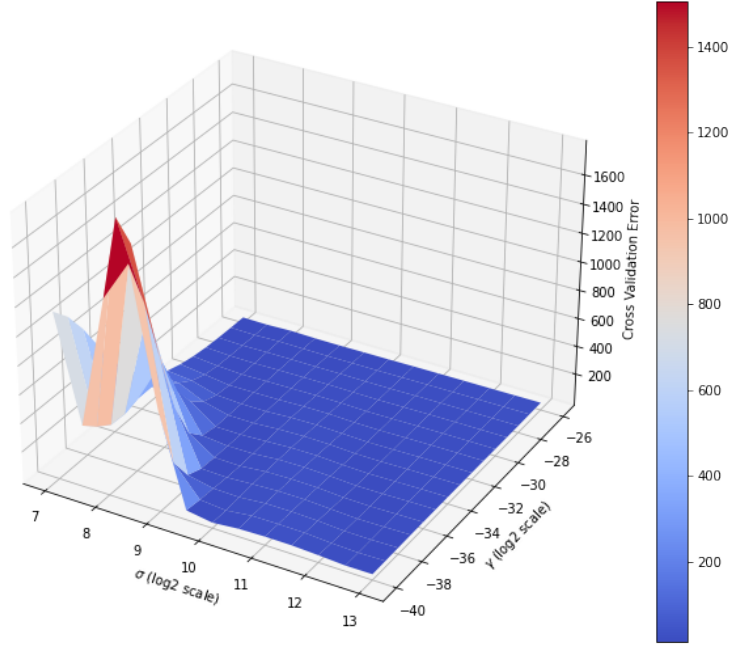


Figure 5: Surface plot of cross validation error as a function of σ and γ values. For convenience, the σ and γ values are plotted on a log-2 scale.

- (b) The cross-validation error as a function of the σ and γ values is plotted as a surface plot in Figure 5.
- (c) The best (σ, γ) pair was defined to be the one that had the least cross validation MSE. This was found to be $\sigma = 2^{9.5}$ and $\gamma = 2^{-36.0}$. Using these values, we again performed kernel ridge regression, this time using the training set and testing set. The train and test MSE for these hyperparameters is as follows:

$$\begin{aligned} \text{Train MSE} &= 5.035038540967276 \\ \text{Test MSE} &= 16.479302486525636 \end{aligned}$$

- (d) Finally, all the different experiments from 4(a,c,d) and 5(a,c) were repeated for 20 random train/test splits to smooth out the results. For convenience, all the experiments for Naive Regression, Linear Regression and Kernel Ridge Regression were converted into functions, which are presented in Appendix 3.2 of the report.

The summarized results are presented in the following table.

Method	MSE Train	MSE Test
Naive Regression	83.517 ± 3.055	86.421 ± 6.143
attribute 1	69.581 ± 2.623	77.042 ± 5.180
attribute 2	72.068 ± 3.240	76.581 ± 6.554
attribute 3	63.221 ± 3.033	67.842 ± 6.142
attribute 4	81.172 ± 3.244	84.269 ± 6.907
attribute 5	67.346 ± 3.129	72.556 ± 6.382
attribute 6	44.215 ± 2.740	42.725 ± 5.437
attribute 7	70.518 ± 3.365	76.685 ± 6.933
attribute 8	77.569 ± 3.298	82.644 ± 6.792
attribute 9	70.002 ± 3.143	76.719 ± 6.460
attribute 10	63.563 ± 3.278	70.825 ± 6.664
attribute 11	61.310 ± 2.796	65.757 ± 5.560
attribute 12	37.638 ± 1.791	40.392 ± 3.651
all attributes	21.736 ± 1.472	25.097 ± 2.926
Kernel Ridge Regression	7.229 ± 1.206	13.738 ± 1.577

Table 4: Summarized results for all regression experiments averaged over 20 runs. Results reported as mean \pm standard deviation.

2 PART II

2.1 Questions

6. Bayes estimator

(a) Let $\hat{y} := f(x)$.

Deriving the optimal solution (Bayes estimator) f^* ,

$$\mathcal{E}(f) := \sum_{x \in X} \sum_{y \in Y} [y \neq f(x)] c_y p(x, y)$$

and $\mathcal{E}(f)$ is the expected error of an arbitrary predictor f .

Separating x and y :

$$\mathcal{E}(f) = \sum_{x \in X} \left[\sum_{y \in Y} [y \neq f(x)] c_y p(y|x) \right] p(x)$$

which can be re-written as:

$$\mathcal{E}(f) = \sum_{x \in X} \left[\sum_{y \neq f(x)} c_y p(y|x) \right] p(x)$$

Now, let us find the value of f^* at a specific point, x' by setting $x := x'$. Then we have:

$$\begin{aligned} \mathcal{E}(f(x')) &= \left[\sum_{y \neq f(x')} c_y p(y|x') \right] p(x') \\ &= \sum_{y \neq f(x')} c_y p(y|x') \quad \text{since } p(x') \text{ is constant} \\ &= 1 - \sum_{y=f(x')} c_y p(y|x') \end{aligned}$$

To find the minimum of \mathcal{E} , we can maximise $c_y p(y|x')$ i.e. $\operatorname{argmax} c_y p(y|x')$.

Therefore,

$$f^* = \operatorname{argmax}_{y=f(x)} c_y p(y = f(x)|x)$$

where the loss is 0 if we predict correctly and c_y otherwise.

(b) i.

$$L_F(y, \hat{y}) = \hat{y}^2 - y^2 + (y - \hat{y})2y = \hat{y}^2 - y^2 + 2y^2 - 2y\hat{y} = y^2 - 2y\hat{y} + \hat{y}^2 = (y - \hat{y})^2$$

This is clearly a square loss function.

ii. Proof of " \Rightarrow ":

if $y = \hat{y}$,

$$L_F(\hat{y}, \hat{y}) = F(\hat{y}) - F(\hat{y}) + (\hat{y} - \hat{y})F'(\hat{y}) = 0 + 0F'(\hat{y}) = 0$$

Hence, we have shown that $y = \hat{y} \rightarrow L_F(y, \hat{y}) = 0$.

Proof of " \Leftarrow ":

if $L_F(y, \hat{y}) = 0$,

$$\begin{aligned} L_F(y, \hat{y}) &= F(\hat{y}) - F(y) + (y - \hat{y})F'(y) = 0 \\ \Rightarrow F(y) &= F(\hat{y}) + (y - \hat{y})F'(y) \end{aligned} \quad (1)$$

Since F is strictly convex,

$$F(\alpha p + (1 - \alpha)q) < \alpha F(p) + (1 - \alpha)F(q) \quad \forall p, q \in \mathbb{R}, \alpha \in (0, 1)$$

$$\begin{aligned} \Rightarrow F(q + \alpha(p - q)) &< F(q) + \alpha(F(p) - F(q)) \\ \Rightarrow F(p) - F(q) &> \frac{F(q + \alpha(p - q)) - F(q)}{\alpha} \end{aligned}$$

As $\alpha \rightarrow 0$, we obtain:

$$\begin{aligned} F(p) - F(q) &> F'(q)(p - q) \\ \Rightarrow F(p) &> F'(q)(p - q) + F(q) \end{aligned} \quad (2)$$

Replacing p, q with y, \hat{y} respectively, we get:

$$F(y) > F'(\hat{y})(y - \hat{y}) + F(\hat{y})$$

Therefore, the only way (1) holds true is if $y = \hat{y}$:

$$F(\hat{y}) = 0F'(\hat{y}) + F(\hat{y}) = F(\hat{y})$$

Hence, we have shown that $L_F(y, \hat{y}) = 0 \rightarrow y = \hat{y}$.

- iii. Replace p, q in equation 2 in 6(b)ii with \hat{y}, y respectively. Then we get:

$$\begin{aligned} F(\hat{y}) &> F'(y)(\hat{y} - y) + F(y) \\ \Rightarrow F(\hat{y}) - F'(y)(\hat{y} - y) - F(y) &> 0 \\ \Rightarrow F(\hat{y}) - F(y) + (y - \hat{y})F'(y) &> 0 \\ \Rightarrow L_F(y, \hat{y}) &> 0 \end{aligned}$$

Thus, we have showed that $L_F(y, \hat{y}) > 0$.

Now, to show $L_F(y, \hat{y}) = 0$, we can use the proof from 6(b)ii that $L_F(y, \hat{y}) = 0$ when $y = \hat{y}$.

Therefore, we have proved that

$$L_F(y, \hat{y}) \geq 0 \quad \forall y, \hat{y} \in \mathbb{R}$$

- iv. Let $\hat{y} := f(x)$.

Deriving the optimal solution (Bayes estimator) f^* in the continuous case,

$$\mathcal{E}(f) := \int_{x \in X} \int_{y \in Y} [F(f(x)) - F(y) + (y - f(x))F'(y)] dP(x, y)$$

and $\mathcal{E}(f)$ is the expected error of an arbitrary predictor f .

Separating x and y :

$$\mathcal{E}(f) = \int_{x \in X} \left[\int_{y \in Y} [F(f(x)) - F(y) + (y - f(x))F'(y)] dP(y|x) \right] dP(x)$$

Now, let us find the value of f^* at a specific point, x' by setting $x := x'$. Then we have:

$$\begin{aligned} \mathcal{E}(f(x')) &= \left[\int_{y \in Y} [F(f(x')) - F(y) + (y - f(x'))F'(y)] dP(y|x') \right] dP(x') \\ &= \int_{y \in Y} [F(f(x')) - F(y) + (y - f(x'))F'(y)] dP(y|x') \end{aligned}$$

since $dP(x')$ is constant

Let $e := \mathcal{E}(f(x'))$ and $z := f(x')$. Replace "=" with " \propto ".

$$e \propto \int_{y \in Y} [F(z) - F(y) + (y - z)F'(y)] dP(y|x')$$

To find the minimum, we differentiate e w.r.t. z :

$$\frac{\partial e}{\partial z} = \int_{y \in Y} [F'(z) - F'(y)] dP(y|x') \quad \text{using derivative rules}$$

Setting it equal to zero and solving,

$$\begin{aligned} 0 &= \int_{y \in Y} F'(z) dP(y|x') - \int_{y \in Y} F'(y) dP(y|x') \\ &= F'(z) - \int_{y \in Y} F'(y) dP(y|x') \\ &\quad \text{since } F'(z) \text{ is independent of } dP(y|x') \text{ and so } \int_{y \in Y} dP(y|x') = 1 \\ \Rightarrow F'(z) &= \int_{y \in Y} F'(y) dP(y|x') \quad (1) \\ \Rightarrow F(z) &= \int_{y \in Y} F'(y) dP(y|x') z + c \quad \text{where } c \text{ denotes constant} \\ \Rightarrow z &= \frac{F(z) - c}{\int_{y \in Y} F'(y) dP(y|x')} \\ \Rightarrow z &= \frac{F(z) - c}{\int_{y \in Y} F'(y) p(y|x') dy} \end{aligned}$$

and by using the conditional distribution function, we can find the Bayes estimator with respect to the the probability mass function $p(x, y)$:

$$f^* = \frac{F(f(x)) - c}{\sum_{y \in Y} F'(y) p(y|x)} \quad \text{for some constant } c$$

In the special case of $F(x) = |x|^p$, using equation 1 and $F'(y) =$

$$p|y|^{p-1}:$$

$$\begin{aligned} p|z|^{p-1} &= \int_{y \in Y} p|y|^{p-1} dP(y|x') \\ \Rightarrow |z|^{p-1} &= \int_{y \in Y} |y|^{p-1} dP(y|x') \\ \Rightarrow |z| &= \left(\int_{y \in Y} |y|^{p-1} dP(y|x') \right)^{\frac{1}{p-1}} \\ \Rightarrow f^* &= \left(\sum_{y \in Y} |y|^{p-1} p(y|x) \right)^{\frac{1}{p-1}} \end{aligned}$$

again by using the conditional distribution function

7. Kernel modification

- (a) By the theorem, K_c is positive semidefinite if and only if $K_c(\mathbf{x}, \mathbf{z}) = \langle \phi(x), \phi(z) \rangle$, where $x, z \in \mathbb{R}^n, c \in \mathbb{R}$.

So assuming that $K_c(\mathbf{x}, \mathbf{z}) = \langle \phi(x), \phi(z) \rangle$, then we have:

$$\begin{aligned} \sum_{i,j=1}^m c_i c_j (c + \sum_{k=1}^n x_{i_k} x_{j_k}) &= c \sum_{i,j} c_i c_j + \sum_{i,j} \sum_k c_i c_j x_{i_k} x_{j_k} \\ &= \sum_k \sum_i c_i x_{i_k} \sum_j c_j x_{j_k} + c \left(\sum_i c_i \sum_j c_j \right) \\ &= \sum_k \left(\sum_i c_i x_{i_k} \right)^2 + c \left(\sum_i c_i \right)^2 \end{aligned}$$

Since $\sum_k (\sum_i c_i x_{i_k})^2 \geq 0$ and $(\sum_i c_i)^2 \geq 0$, for the entire equation to be ≥ 0 , c must be ≥ 0 .

Hence, K_c is a positive semidefinite kernel for $c \geq 0$. The proof for this has already been done above.

- (b) c in K_c acts like a bias term, b in linear regression equation, $\mathbf{y} = X\mathbf{w} + b$. It allows us to get the best fit offset from the origin. Without it, we are forcing the line to go through the origin which might not provide the best fit and bias other parameters.

8. Under a scenario in which m , a number of data is small, we select β to be a function of $\mathbf{x}_i \in \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$, and the test point \mathbf{t} :

$$\beta = \hat{\beta}(\mathbf{x}_i, \mathbf{t}) = \begin{cases} \infty & \text{if } \mathbf{x}_i \text{ is the nearest point to } \mathbf{t} \\ -\infty & \text{otherwise} \end{cases}$$

Here, the nearest point \mathbf{x}_i to \mathbf{t} can be found by calculating the distance between \mathbf{x}_i and \mathbf{t} i.e. $\|\mathbf{x}_i - \mathbf{t}\|$.

Argument:

One of differences between this trained linear classifier and 1-NN classifier is that the former considers all the data points whereas the latter considers one data point that is closest to the test point. To make the trained linear classifier behave like a 1-Nearest Neighbor classifier i.e. to make the trained linear classifier to take only one nearest data point to the test point into account, we want only one nearest neighbour to receive the highest weight and the rest to receive the lowest possible weight so that they would almost be 'ignored'. This is where the Gaussian kernel comes into play to act as that weight. Since the Gaussian kernel output depends on β , we can define β to return a value such that the kernel function output is the largest with an input, a closest data point to the test point and smallest with other inputs. Since the kernel function decays exponentially, we know that $\exp(-\beta\|\mathbf{x} - \mathbf{t}\|^2) \rightarrow \infty$ as $-\beta\|\mathbf{x} - \mathbf{t}\|^2 \rightarrow -\infty$ and $\exp(-\beta\|\mathbf{x} - \mathbf{t}\|^2) \rightarrow 0$ as $-\beta\|\mathbf{x} - \mathbf{t}\|^2 \rightarrow \infty$ which means that when \mathbf{x} is the nearest neighbour to \mathbf{t} , β needs to return ∞ and $-\infty$ otherwise. Now, if we assume \mathbf{x}_c is a nearest point to \mathbf{t} then we would have $f(\mathbf{t}) \approx \alpha_c K_\beta(\mathbf{x}_c, \mathbf{t})$ for $c \in \{1, \dots, m\}$ since $\alpha_i K_\beta(\mathbf{x}_c, \mathbf{t}) \approx 0$ for all $i \neq c$. So $\text{sign}(f(\mathbf{t}))$ would be equal to y_c .

9. Let us represent:

- a $n \times n$ matrix, C , as a board configuration indicating non-empty holes as 1s and empty holes as 0s. This matrix will be initialised to an initial board configuration, G .
- a holes hit as a $n \times n$ matrix, A , representing holes being hit in row i and column j . For example, given $n = 4$:

$$A_{23} = \begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

Given above, we can represent a $n \times n$ matrix, C' , showing a new board configuration after one hit, A_{ij} to the previous board before A_{ij} , C :

$$C' = A_{ij} + C$$

where each entry is in modulo 2. This means that $1 + 1 = 2 = 0 \pmod{2}$, which is correct as any moles currently appearing will hide after a hit. Similarly, $1 + 0 = 1 = 1 \pmod{2}$ is also correct as new moles pop up when empty holes are hit.

This means that if a sequence of holes, that you can hit to empty the board, exists then there is a combination of matrices, when added to the initial board configuration, G , that produces a zero matrix representing an empty board:

$$G + \sum_{i,j=1}^n \alpha_{ij} A_{ij} = 0 \tag{1}$$

where $\alpha_{ij} \in \{0, 1\}$ is a coefficient deciding whether to hit the hole or not. Therefore, the task would be to find such a sequence of α_{ij} to determine whether the solution exists or not. In other words, a solution exists if there is a sequence of hit patterns to whack all the moles in G . Re-arranging the equation (1) above, we obtain:

$$G = - \sum_{i,j=1}^n \alpha_{ij} A_{ij} = \sum_{i,j=1}^n \alpha_{ij} A_{ij} \quad \text{since } -1 = 1 \pmod{2} \tag{2}$$

which can be expanded to:

$$G = \alpha_{11} \begin{pmatrix} 1 & 1 & 0 & \dots & 0 \\ 1 & 0 & 0 & \dots & 0 \\ 0 & & \dots & & 0 \\ \vdots & & \ddots & & \vdots \\ 0 & & \dots & & 0 \end{pmatrix} + \alpha_{12} \begin{pmatrix} 1 & 1 & 1 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & & \dots & & 0 \\ \vdots & & \ddots & & \vdots \\ 0 & & \dots & & 0 \end{pmatrix} + \dots$$

$$+ \alpha_{nn} \begin{pmatrix} 0 & & \dots & & 0 \\ \vdots & & \ddots & & \vdots \\ 0 & & \dots & & 0 \\ 0 & \dots & 0 & 0 & 1 \\ 0 & \dots & 0 & 1 & 1 \end{pmatrix}$$

Here, we can notice that each entry in G can be represented as a linear equation. For example:

$$G_{22} = \alpha_{11} + \alpha_{21} + \dots + \alpha_{32}$$

This tells us that we can represent G as a $n^2 \times 1$ vector of its entries as follows:

$$G = \begin{bmatrix} G_{11} \\ G_{12} \\ \vdots \\ G_{nn} \end{bmatrix}$$

Thus, the equation (2) can be re-written as a system of linear equations, $\mathbf{A}\vec{\alpha} = \vec{G}$:

$$[vec(A_{11})vec(A_{12}) \dots vec(A_{nn})] \cdot \begin{bmatrix} \alpha_{11} \\ \alpha_{12} \\ \vdots \\ \alpha_{nn} \end{bmatrix} = \begin{bmatrix} G_{11} \\ G_{12} \\ \vdots \\ G_{nn} \end{bmatrix}$$

where $vec(A_{ij})$ denotes a vectorised matrix A_{ij}

To find $\vec{\alpha}$, we can use [Gaussian Elimination](#) (row reduction). Given $\vec{\alpha}$, to obtain the sequence of hit patterns, we can find each solution's (α_{ij}) corresponding A_{ij} where $\alpha_{ij} = 1$. Convert each A_{ij} to an integer labelling each hole on the board as shown in the question and we are done.

We can show that the algorithm is polynomial in n :

First, computing each A_{ij} in \mathbf{A} takes $O(n^2)$ and since there are n^2 number of A_{ij} s, overall, computing \mathbf{A} takes $O(n^4)$. Performing the Gaussian elimination takes [O\(\$n^6\$ \)](#), which dominates the runtime of the entire algorithm. Hence, we have showed that the algorithm is polynomial in n .

3 Appendices

3.1 Vectorized Implementation of Kernel Ridge Regression

Here, we present the mathematical formulation for the vectorized implementation of kernel ridge regression. This allows us to avoid using for-loops, thus improving efficiency of training.

3.1.1 Kernel Matrix

We first vectorize the computation of the kernel matrix \mathbf{K} , which is required to compute the dual solution α^* . To do this, we begin with the original description of \mathbf{K} :

$$\mathbf{K} = \begin{bmatrix} K(x_1, x_1) & \cdots & K(x_1, x_\ell) \\ \vdots & \ddots & \vdots \\ K(x_\ell, x_1) & \cdots & K(x_\ell, x_\ell) \end{bmatrix}$$

$$K(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right) = \exp\left(-\frac{\langle x_i - x_j, x_i - x_j \rangle}{2\sigma^2}\right)$$

Note that ℓ is the number of training examples, thus \mathbf{K} is an $(\ell \times \ell)$ matrix. Using the NumPy framework, we can reformulate \mathbf{K} as:

$$\mathbf{K} = \text{np.exp}\left(-\frac{1}{2\sigma^2} \begin{bmatrix} \langle x_1 - x_1, x_1 - x_1 \rangle & \cdots & \langle x_1 - x_\ell, x_1 - x_\ell \rangle \\ \vdots & \ddots & \vdots \\ \langle x_\ell - x_1, x_\ell - x_1 \rangle & \cdots & \langle x_\ell - x_\ell, x_\ell - x_\ell \rangle \end{bmatrix}\right)$$

The goal now is to breakdown the matrix in the above expression into a combination of simpler matrices. Note that a general inner product in this matrix can be expanded as follows:

$$\langle x_i - x_j, x_i - x_j \rangle = \langle x_i, x_i \rangle - 2\langle x_i, x_j \rangle + \langle x_j, x_j \rangle \quad (5)$$

We can use this property to redefine the matrix as follows:

$$\begin{bmatrix} \langle x_1 - x_1, x_1 - x_1 \rangle & \cdots & \langle x_1 - x_\ell, x_1 - x_\ell \rangle \\ \vdots & \ddots & \vdots \\ \langle x_\ell - x_1, x_\ell - x_1 \rangle & \cdots & \langle x_\ell - x_\ell, x_\ell - x_\ell \rangle \end{bmatrix} = \mathbf{B} - 2\mathbf{X}\mathbf{X}^T + \mathbf{B}^T \quad (6)$$

where

$$\mathbf{B} = \begin{bmatrix} \langle x_1, x_1 \rangle & \cdots & \langle x_1, x_\ell \rangle \\ \vdots & \ddots & \vdots \\ \langle x_\ell, x_1 \rangle & \cdots & \langle x_\ell, x_\ell \rangle \end{bmatrix}, \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_\ell^T \end{bmatrix} \quad (7)$$

Essentially what we have done is break down each element of the LHS of 6 into the 3 terms as per equation 5 and then the first, second and third terms are separated into individual matrices. The second matrix consists of the terms $-2\langle x_i - x_j, x_i - x_j \rangle$, which can be simplified to the matrix $-2\mathbf{X}\mathbf{X}^T$ where \mathbf{X} is just the $(l \times n)$ training data matrix. It can also be noticed that the first and third matrices are just transposes (hence \mathbf{B} and \mathbf{B}^T) of each other and contain the diagonal terms of the $\mathbf{X}\mathbf{X}^T$ matrix stacked ℓ times as rows or columns. When implementing in NumPy however, we could just represent \mathbf{B} as a single column of the diagonals of $\mathbf{X}\mathbf{X}^T$. Due to broadcasting, we can expect that the vector and matrix additions will still work. Thus, we summarize the \mathbf{K} reformulation as

$$\mathbf{K} = \text{np.exp} \left(-\frac{1}{2\sigma^2} (\mathbf{B} - 2\mathbf{X}\mathbf{X}^T + \mathbf{B}^T) \right)$$

3.1.2 Evaluation function

The next component is vectorizing the evaluation function. This function has the form

$$y_{test} = \sum_{i=1}^l \alpha_i^* K(x_i, x_{test})$$

where y_{test} is the predicted value for a test example x_{test} . Using the same logic as in the kernel matrix reformulation, we get

$$y_{test} = \mathbf{K}_{eval} \alpha^*$$

where

$$\mathbf{K}_{eval} = \text{np.exp} \left(-\frac{1}{2\sigma^2} (\mathbf{B}^T - 2\mathbf{X}_{test}\mathbf{X}^T + \mathbf{C}) \right)$$

where \mathbf{X}_{test} is an $(m \times n)$ testing data matrix, \mathbf{B}^T is same as defined in 7 and \mathbf{C} is defined as the diagonals of matrix $\mathbf{X}_{test}\mathbf{X}_{test}^T$. Note however that the \mathbf{K}_{eval} matrix is not symmetric matrix like the kernel matrix; it is of dimensions $(m \times l)$.

3.2 Functions for regression experiments on Boston Housing Dataset

```
1 def naive_regression_eval(num_train_samples, num_test_samples
2   , train_y, test_y, results_dict):
3     """
4     Performs Naive Regression on given data and saves results.
5
6     Args
7     ----
8     num_train_samples : number of training samples. This will
9       be used to create a ones vector for the training set.
10    num_test_samples : number of testing samples. This will be
11      used to create a ones vector for the testing set.
12    train_y : NumPy array of training output values. Should be
13      of shape (num_train_samples, 1)
14    test_y : NumPy array of testing output values. Should be of
15      shape (num_test_samples, 1)
16    results_dict : Dictionary in which to save the results.
17
18    Returns
19    -----
20    None.
21
22    """
23
24    #Add the key into the results dictionary
25    key = "Naive Regression"
26    if not key in results_dict.keys():
27        results_dict[key] = {"Train MSE": [], "Test MSE": []}
28
29    #Create a ones vector for train and test data
30    train_ones = np.ones((num_train_samples, 1))
31    test_ones = np.ones((num_test_samples, 1))
```



```

27
28 #Compute the weights using training set - in case we
    encounter singular matrices, Ive added the psuedoinverse
    as backup
29 try:
30     w = np.linalg.inv(train_ones.T @ train_ones) @ train_ones
        .T @ train_y
31 except np.linalg.LinAlgError:
32     w = np.linalg.pinv(train_ones.T @ train_ones) @
        train_ones.T @ train_y
33
34 #Add results into the dictionary
35 results_dict[key]["Train MSE"].append(MSE(train_ones,
        train_y, w).item())
36 results_dict[key]["Test MSE"].append(MSE(test_ones, test_y,
        w).item())
37
38 print(">> Naive Regression done.")
39
40 def linear_regression_eval(train_x, test_x, train_y, test_y,
    results_dict, single=False):
41     """
42     Performs linear regression on the given data and saves
        results.
43
44     Args
45     ----
46     train_x : NumPy array of training inputs. Should be of
        shape (m, n) where m is number of training examples and n
        is number of input features
47     test_x : NumPy array of test inputs. Should be of shape (l,
        n) where l is number of testing examples and n is number
        of input features
48     train_y : NumPy array of training outputs. Should be of
        shape (m, 1) where m is number of training examples
49     test_y : NumPy array of testing inputs. Should be of shape
        (l, 1) where l is number of training examples
50     results_dict : Dictionary in which to save the results.
51     single : Set to True if you wish to perform on single
        attributes.
52
53     Returns
54     -----
55     None.
56     """

```

```

57
58 if single:
59     num_features = train_x.shape[1] - 1
60
61     #Iterate over the features
62     for i in range(num_features):
63         #Create a key for the attribute and add it into the
        results dictionary if it doesn't exist
64         key = "attribute {}".format(i+1)
65         if not key in results_dict.keys():
66             results_dict[key] = {"Train MSE": [], "Test MSE": []}
67
68         #Pick out a the specific feature and the bias column
69         train_x_subset = train_x[:, [0, i+1]]
70         test_x_subset = test_x[:, [0, i+1]]
71
72         #Checkpoint to make sure the dimensions are correct
73         assert train_x_subset.shape == (train_x.shape[0], 2)
74
75         #Computed least squares solution - In case we encounter
        singular matrices, I have used the pseudoinverse as
        backup
76         try:
77             w = np.linalg.inv(train_x_subset.T @ train_x_subset)
78             @ train_x_subset.T @ train_y
79             except np.linalg.LinAlgError:
80                 w = np.linalg.pinv(train_x_subset.T @ train_x_subset)
81                 @ train_x_subset.T @ train_y
82
83         # Compute the train and test MSE and add it into the
        results_dict
84         results_dict[key]["Train MSE"].append(MSE(
85         train_x_subset, train_y, w).item())
86         results_dict[key]["Test MSE"].append(MSE(test_x_subset,
87         test_y, w).item())
88
89         print(">> Linear regression for {} done.".format(key))
90
91 else:
92     #Create a key for the attribute and add it into the
        results dictionary if it doesn't exist
93     key = "all attributes"
94     if not key in results_dict.keys():
95         results_dict[key] = {"Train MSE": [], "Test MSE": []}
96

```

```

93     #Compute the weights using training set - in case we
    encounter singular matrices, Ive added the psuedoinverse
    as backup
94     try:
95         w = np.linalg.inv(train_x.T @ train_x) @ train_x.T @
    train_y
96     except np.linalg.LinAlgError:
97         w = np.linalg.pinv(train_x.T @ train_x) @ train_x.T @
    train_y
98
99     #Compute the train and test MSE and add it into the
    results dictionary
100     results_dict[key]["Train MSE"].append(MSE(train_x,
    train_y, w).item())
101     results_dict[key]["Test MSE"].append(MSE(test_x, test_y,
    w).item())
102
103     print(">> Linear Regression for {} done.".format(key))
104
105 def ridge_regression_eval(train_x, train_y, test_x, test_y,
    sigma_values, gamma_values, results_dict):
106     """
107     Performs kernelized ridge regression on the given data and
    saves results.
108
109     Args
110     ----
111     train_x : NumPy array of training inputs. Should be of
    shape (m, n) where m is number of training examples and n
    is number of input features
112     test_x : NumPy array of test inputs. Should be of shape (l,
    n) where l is number of testing examples and n is number
    of input features
113     train_y : NumPy array of training outputs. Should be of
    shape (m, 1) where m is number of training examples
114     test_y : NumPy array of testing inputs. Should be of shape
    (l, 1) where l is number of training examples
115     sigma_values : NumPy 1D array of possible sigma values.
116     gamma_values : NumPY 1D array of possible gamma values.
117
118     Returns
119     -----
120     None.
121     """
122

```

```

123 #Generate a seed value which we will use inside the for
    loop.
124 seed=np.random.randint(100)
125
126 #Create a key for the attribute and add it into the results
    dictionary if it doesn't exist
127 key = "Kernel Ridge Regression"
128 if not key in results_dict.keys():
129     results_dict[key] = {"Train MSE": [], "Test MSE": []}
130
131 #The best parameters will be saved in the order [best_sigma
    , best_gamma]
132 best_parameters = {"sigma":0, "gamma":0}
133
134 #Initialize the smallest error to be some large number.
    This will be used to update our best parameters
135 min_error = 10000000000000000000
136
137 #Iterate over entire parameter space
138 for i in range(len(sigma_values)):
139     for j in range(len(gamma_values)):
140
141         #Create the sigma and gamma pair
142         sigma = sigma_values[i]
143         gamma = gamma_values[j]
144
145         #Perform 5-fold cross validation.
146         # Note - The train_with_kfoldCV() function performs a
            kfoldCV for just one sigma and gamma.
147         # In the function, I've added a shuffle operation
            BEFORE splitting the data into the k groups.
148         # This means that everytime we use a new sigma and
            gamma, the group compositions may differ due to the
            shuffle.
149         # We need to ensure that the group compositions after
            shuffling are the same everytime so that our parameters
            sigma and gamma are comparable.
150         np.random.seed(seed)
151         mean_cv_error = train_with_kfoldCV(k=5, data=(train_x,
            train_y), kernel_params=(sigma, gamma), shuffle=True,
            verbose=False)
152
153         #If the mean_cv_error is smaller than the minimum error
            , then we can update our best parameters
154         if mean_cv_error < min_error:

```

```

155         min_error = mean_cv_error
156         best_parameters['sigma'] = sigma
157         best_parameters['gamma'] = gamma
158
159     #Compute dual regression solution using training set
160     alpha_star = kernel_ridge_regression(X=train_x, y=train_y,
161                                         sigma=best_parameters['sigma'], gamma=best_parameters['
162                                         gamma'])
163
164     #Compute predictions on train set: The "test" points here
165     #are just the training set again
166     train_y_pred = evaluate(alpha=alpha_star, X_tr=train_x,
167                             X_te=train_x, sigma=best_parameters['sigma'])
168     #Compute predictions on test set
169     test_y_pred = evaluate(alpha=alpha_star, X_tr=train_x, X_te
170                             =test_x, sigma=best_parameters['sigma'])
171
172     #Compute train and test MSE
173     train_MSE = (1/len(train_y)) * (train_y_pred - train_y).T @
174                 (train_y_pred - train_y)
175     test_MSE = (1/len(test_y)) * (test_y_pred - test_y).T @ (
176                 test_y_pred - test_y)
177
178     results_dict[key]["Train MSE"].append(train_MSE)
179     results_dict[key]["Test MSE"].append(test_MSE)
180
181     print(">> Kernel Ridge Regression done.\n")
182
183     #Initialize a final results dictionary which we shall update
184     #with all the errors
185     final_results = dict()
186
187     #Store number of runs here
188     runs = 20
189
190     #Generate a seed for reproducible results
191     np.random.seed(1753295)
192
193     #PROBLEM: INSPITE OF SETTING SEED HERE, THE FUNCTION IS
194     #GENERATING SAME TRAINING SET EVERYTIME!! HOW TO FIX IT?
195
196     #Iterate over runs
197     for i in range(runs):
198
199         print("Run {}".format(i+1))

```

```

191
192 #Generate a fresh batch of train/test data
193 train_X, train_Y, test_X, test_Y = split_data(inputs=X,
194         targets=Y, test_proportion=1/3, shuffle=True)
195 # print("Training set = ", train_X)
196
197 #PERFORM NAIVE REGRESSION AND UPDATE RESULTS
198 naive_regression_eval(num_train_samples=len(train_X),
199         num_test_samples=len(test_X), train_y=train_Y, test_y=
200         test_Y, results_dict=final_results)
201
202 #PERFORM LINEAR REGRESSION ON SINGLE ATTRIBUTES AND UPDATE
203         RESULTS
204 linear_regression_eval(train_x=train_X, test_x=test_X,
205         train_y=train_Y, test_y=test_Y, results_dict=final_results
206         , single=True)
207
208 #PERFORM LINEAR REGRESSION ON ALL ATTRIBUTES AND UPDATE
209         RESULTS
210 linear_regression_eval(train_x=train_X, test_x=test_X,
211         train_y=train_Y, test_y=test_Y, results_dict=final_results
212         , single=False)
213
214 #PERFORM KERNELIZED RIDGE REGRESSION ON ALL ATTRIBUTES AND
215         UPDATE RESULTS
216 #Create parameter space
217 sigma_values = 2*np.arange(7, 13.5, 0.5)
218 gamma_values = (2*np.arange(40.0, 25.0, -1))**-1
219 # sigma_values = 2*np.arange(7, 8.0, 0.5)
220 # gamma_values = (2*np.arange(40.0, 38.0, -1))**-1
221 ridge_regression_eval(train_x=train_X[:, 1:], train_y=
222         train_Y,
223         test_x=test_X[:, 1:], test_y=test_Y,
224         sigma_values=sigma_values,
225         gamma_values=gamma_values,
226         results_dict=final_results)
227 print("Evaluation Completed.")

```