

Supervised Learning (COMP0078) – Coursework 1

Bowen Ding (21118015)
Zhengxiang Wei (21181074)

November 15, 2021

Abstract

This is the answer report for Supervised Learning (COMP0078) Coursework 1. Answers are arranged in the order of Part I, Part II and Part III. All the code are attached at the end of this report with a single jupyter notebook file attached in the submission of work.

Contents

1	Part I	2
1.1	Linear Regression	2
1.2	Filtered Boston housing and kernels	5
1.3	Kernelised ridge regression	6
2	Part II	8
2.1	k -Nearest Neighbors	8
2.1.1	Generating the data	8
2.1.2	Estimated generalization error of k -NN as a function of k	8
2.1.3	Determine the optimal k -NN as a function of the number of training points (m)	9
2.1.4	Additional Clarifications	10
3	Part III	10
3.1	Questions	10
A	Appendix I: Code for Part I	12
A.1	Plot fit function curve, give corresponding equation and calculate mse for Q1	12
A.2	Plot mse(or Ln(mse)) for both training and test set in Q2 to explore overfitting	13
A.3	Plot mse(or $\ln(\text{mse})$) for both training and test set based on a different basis function for Q3	16
A.4	Baseline versus full linear regression	17
A.5	Perform KRR ,find the best pair of gamma and sigma for KRR, compare the performance with other methods.	19
B	Appendix II: Code for Part II	23
B.1	K-NN Classifier and generating points	23
B.2	Do k-nn labelling and plot question 6 figure	24
B.3	Do k-nn test and plot question 7 figure	24
B.4	Do k-nn test and plot question 8 figure	25

Introduction

In this submission, our coding is done with jupyter notebook, using python 3.9 as the kernel language. As required, we have implemented the Regression and $k - NN$ by writing our own functions and classes. To clarify that no packages that was not allowed were used in our work, here is a list of packages we used in this coursework:

- Numpy, for mathematical processing with data such as calculating, sorting and etc.;
- Numba.njit, for accelerating the speed of running the code in part II with lots for loops;
- matplotlib.pyplot, to give plots required in the coursework;

1 Part I

In this section question 1-5 would be answered.

1.1 Linear Regression

Considering a known data set of $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$, where $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{in})^T$ is a input sample point with n dimensions while corresponding y_i is a real number, which is the label of \mathbf{x}_i . It is shown in data set D , m points with labels are sampled independently and from an identical distribution. In linear regression with basis functions, with the help of basis functions $(\phi_1(\mathbf{x}), \dots, \phi_N(\mathbf{x}))$, where $\phi_i(\mathbf{x}) : R^n \rightarrow R$, we can define a feature map $\phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_k(\mathbf{x}))^T$ and mapped data matrix $A \in R^{m \times k}$:

$$A := \begin{bmatrix} \phi(\mathbf{x}_1)^T \\ \vdots \\ \phi(\mathbf{x}_m)^T \end{bmatrix} = \begin{bmatrix} \phi_1(\mathbf{x}_1) & \dots & \phi_k(\mathbf{x}_1) \\ \vdots & \dots & \vdots \\ \phi_1(\mathbf{x}_m) & \dots & \phi_k(\mathbf{x}_m) \end{bmatrix}$$

The main idea of linear regression is to find function $f(\mathbf{x}_i) = w^T \phi(\mathbf{x}_i)$ which is a linear combination of $\phi_j(\mathbf{x}_i)$, $j = 1 \dots k$ to fit y_i . We can calculate regression coefficients by the formula below:

$$w = (A^T A)^{-1} A^T y$$

where $y = (y_1, y_2, \dots, y_m)^T$

Question 1

In this question, we have the data set $\{(1, 3), (2, 2), (3, 0), (4, 5)\}$ after applying k -order polynomial basis functions, we get mapped data matrix A :

$$A = \begin{bmatrix} \mathbf{x}_1^0 & \mathbf{x}_1^1 & \dots & \mathbf{x}_1^j & \dots & \mathbf{x}_1^{k-1} \\ \mathbf{x}_2^0 & \mathbf{x}_2^1 & \dots & \mathbf{x}_2^j & \dots & \mathbf{x}_2^{k-1} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \mathbf{x}_i^0 & \mathbf{x}_i^1 & \dots & \mathbf{x}_i^j & \dots & \mathbf{x}_i^{k-1} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \mathbf{x}_m^0 & \mathbf{x}_m^1 & \dots & \mathbf{x}_m^j & \dots & \mathbf{x}_m^{k-1} \end{bmatrix}$$

(a) We superimpose the curves with polynomial basis of $k=1, 2, 3, 4$ respectively to fit over four data points.

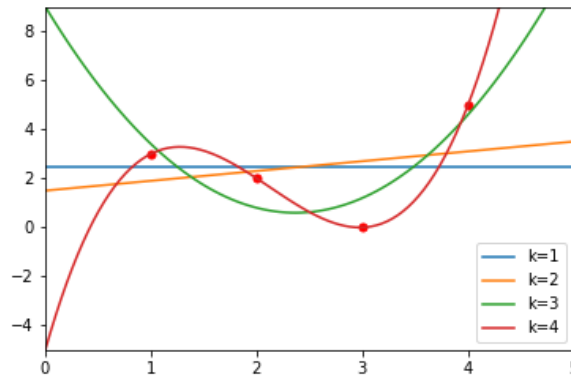


Figure 1: Four different fitting curves with polynomial basis of order=1,2,3,4

(b) The equations corresponding to the curves fitted for $k=1,2,3$ are as below:

- $k=1, f(x) = 2.5$
- $k=2, f(x) = 1.5 + 0.4x$
- $k=3, f(x) = 9 - 7.1x + 1.5x^2$

(c) We use $mse = \frac{1}{m} \sum_{i=1}^m \left(y_i - \sum_{j=1}^k w_j \phi_j(x_i) \right)^2$ to measure loss. The lower the value of mse, the better the algorithm fits.

- $k=1, mse=3.25$
- $k=2, mse=3.05$
- $k=3, mse=0.80$
- $k=4, mse=5.29 \times 10^{-24}$

Question 2

In question2, we explore the overfitting phenomena. Firstly, we generate sample \mathbf{x}_i uniformly at random from the interval $[0, 1]$ 30 times and apply $g_\sigma(\mathbf{x}_i) := \sin^2(2\pi\mathbf{x}_i) + \varepsilon$ to generate training set, where $\varepsilon \sim N(0, \sigma^2)$ serves as noise. The training set is as below:

$$S_{0.07,30} = \{(\mathbf{x}_1, g_{0.07}(\mathbf{x}_1)), (\mathbf{x}_2, g_{0.07}(\mathbf{x}_2)), \dots, (\mathbf{x}_{30}, g_{0.07}(\mathbf{x}_{30}))\}$$

(a) i. The graph of the function $\sin^2(2\pi x)$ and scatter of samples in training set.

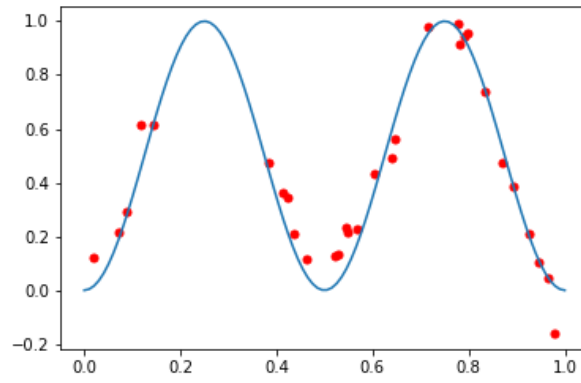


Figure 2: $\sin^2(2\pi x)$ and training set samples

(a) ii. We fit the training set with a polynomial basis of dimension $k=2,5,10,14,18$

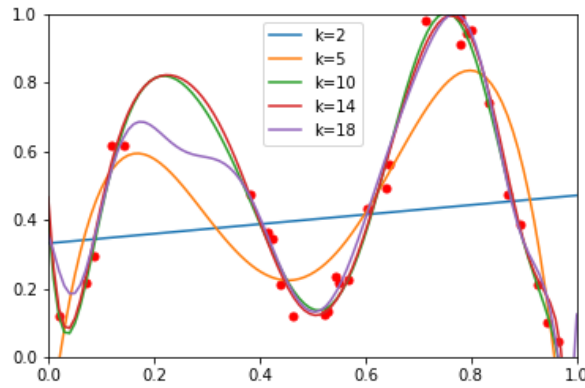


Figure 3: Polynomial fit curves of $k=2,5,10,14,18$

(b) Using $\ln(mse)$ to measure loss in the training set against polynomial dimensions.

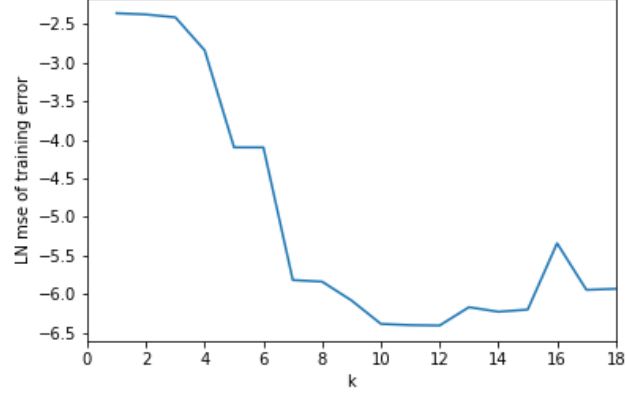


Figure 4: Training set loss against k=1...18

(c) In order to test whether the model begins to fit the noise rather than fitting the function, with model's complexity (dimension of polynomial) increasing, we generate a test set T with 1000 points:

$$T_{0.07,1000} = \{(\mathbf{x}_1, g_{0.07}(\mathbf{x}_1)), (\mathbf{x}_2, g_{0.07}(\mathbf{x}_2)), \dots, (\mathbf{x}_{1000}, g_{0.07}(\mathbf{x}_{1000}))\}$$

Then we use $\ln(mse)$ to measure loss in the test set against polynomial dimensions.

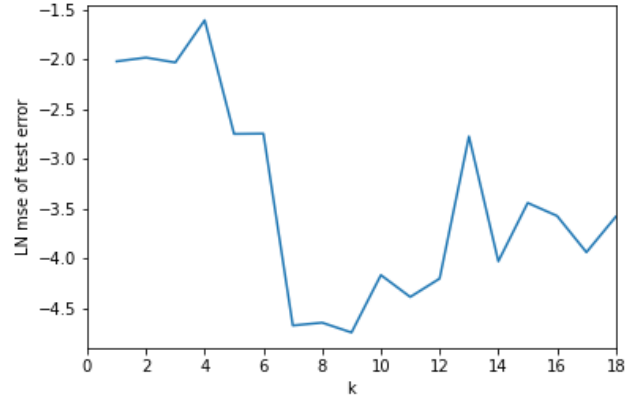


Figure 5: Test set loss against k=1...18

(d) We repeat (b) and (c) 100 times and use $\ln(\text{average results of a 100 runs})$ to measure the loss both in training set and test set.

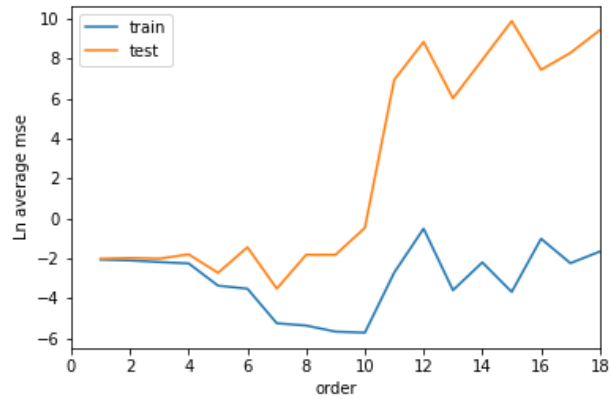


Figure 6: $\ln(\text{average mse of a 100 runs})$ for training and test set

Question 3

In this case, the mapped data matrix A could be:

$$A = \begin{bmatrix} \sin(\pi x_1) & \sin(2\pi x_1) & \cdots & \sin(j\pi x_1) & \cdots & \sin(k\pi x_1) \\ \sin(\pi x_2) & \sin(2\pi x_2) & \cdots & \sin(j\pi x_2) & \cdots & \sin(k\pi x_2) \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \sin(\pi x_i) & \sin(2\pi x_i) & \cdots & \sin(j\pi x_i) & \cdots & \sin(k\pi x_i) \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \sin(\pi x_m) & \sin(2\pi x_m) & \cdots & \sin(j\pi x_m) & \cdots & \sin(k\pi x_m) \end{bmatrix}$$

From the picture below, we can see $\ln(mse)$ to measure loss in the training set and test set against polynomial dimensions for one round.

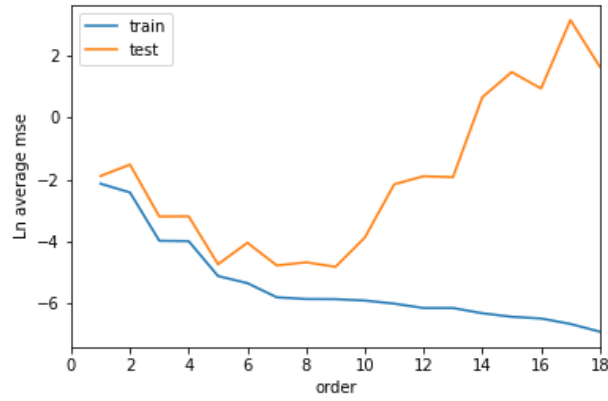


Figure 7: $\ln(mse)$ for training and test set, $\sin(k\pi x)$ basis functions

We also repeat 100 rounds,

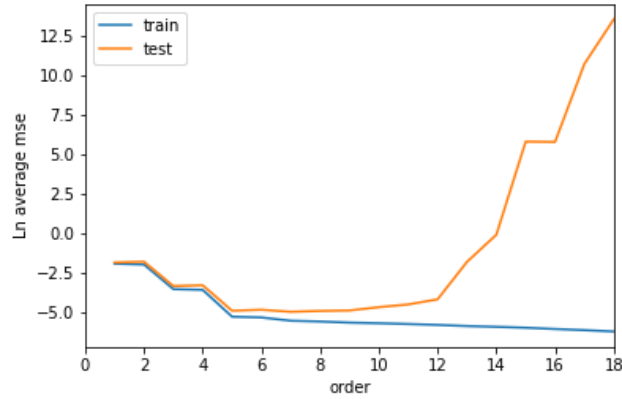


Figure 8: $\ln(\text{average mse of 100 runs})$ for training and test set, $\sin(k\pi x)$ basis functions

1.2 Filtered Boston housing and kernels

Question 4

Compare the performance among different regression methods. We split the data set into training set (2/3 of total data set) and test set (1/3 of total data set), and perform regression over 20 runs and calculate the average mse of every runs.

- Perform naive regression on the training set and calculate the mse on the training and test sets.

In the case of Naive Regression. The basis function of this case is $\phi(x) = \{1\}$, which means mapped matrix A

is :

$$A = \begin{bmatrix} 1 \\ 1 \\ \dots \\ 1 \end{bmatrix} \in R^m$$

- MSE on the training set: 82.87
- MSE on the test set: 86.92

b. Give a simple interpretation of the constant function in a.

Since,

$$w = \left(A^T A\right)^{-1} A^T y$$

Hence,

$$w = \frac{1}{m} \sum_{i=1}^m y_i$$

$$predict = Aw = \begin{bmatrix} \frac{1}{m} \sum_{i=1}^m y_i \\ \frac{1}{m} \sum_{i=1}^m y_i \\ \dots \\ \frac{1}{m} \sum_{i=1}^m y_i \end{bmatrix}$$

We can find that the constant function in (a) is

$$f(\mathbf{x}_i) = \frac{1}{m} \sum_{i=1}^m y_i$$

which is the mean of the label of inputted samples.

c. Linear Regression with single attributes.

The fit function of this case could be

$$f(\mathbf{x}_i) = w_0 + w_1 x_{ij}, i = 1 \dots m, j = 1 \dots 12$$

where x_{ij} means the value of i-th sample point's j-th attribute.

- attribute1: MSE train 70.53869529 MSE test 73.00817819
- attribute2: MSE train 71.96287247 MSE test 75.48257325
- attribute3: MSE train 63.25866244 MSE test 66.65348237
- attribute4: MSE train 80.61347746 MSE test 82.74184527
- attribute5: MSE train 67.55052457 MSE test 70.96812434
- attribute6: MSE train 43.44303058 MSE test 43.13931893
- attribute7: MSE train 70.54381803 MSE test 75.17310599
- attribute8: MSE train 77.51229332 MSE test 81.32391081
- attribute9: MSE train 70.9495549 MSE test 73.43982526
- attribute10: MSE train 64.83087928 MSE test 66.95009183
- attribute11: MSE train 61.58601723 MSE test 63.76119662
- attribute12: MSE train 37.28030312 MSE test 40.30809575

d. Linear Regression using all attributes.

- MSE train: 21.631067081403906
- MSE test: 21.791834593295146

1.3 Kernelised ridge regression

Considering that we have l examples, each sample point has n attributes, so we can try to find a weight of ridge regression by optimisation:

$$w^* = \arg \min_{w \in R^n} \frac{1}{l} \sum_{i=1}^l \left(x_i^T w - y_i \right)^2 + \gamma w^T w$$

For a given kernel function K define the kernel matrix \mathbf{K} :

$$K_{i,j} := K(\mathbf{x}_i, \mathbf{x}_j), i, j = 1 \dots l$$

The dual optimisation formulation after kernelization is

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

We may solve α^* by:

$$\alpha^* = (K + \gamma I_l)^{-1} y$$

and predict y by:

$$y_{test} = \sum_{i=1}^l \alpha_i^* K(\mathbf{x}_i, \mathbf{x}_{test})$$

Question 5

In this exercise, we perform kernel ridge regression on the data set with Gaussian kernel,

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right)$$

We still hold 2/3 of total data for training and 1/3 for testing.

a. We perform kernel ridge regression on the training set using five-fold cross-validation to find the best pair of γ and σ ($\gamma \in [2^{-40}, 2^{-39}, \dots, 2^{-26}]$, $\sigma \in [2^7, 2^{7.5}, \dots, 2^{12.5}, 2^{13}]$), whose test error of five-fold cross-validation is the lowest.

- best gamma: 2^{-33}
- best sigma: 2^{10}

b. We plot the "five-fold cross validation error" as a function of γ and σ

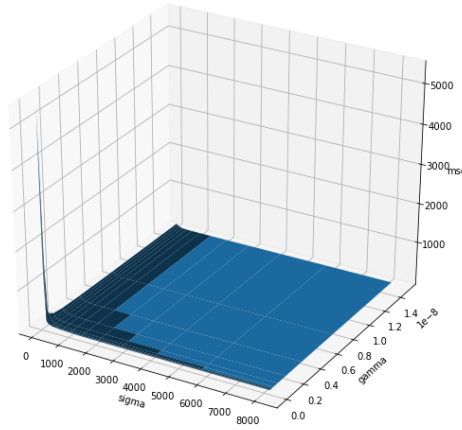


Figure 9: 5-fold cross validation mse of each pair of γ and σ

c. The MSE on training and test sets for the best γ and σ :

- MSE train: 22.634913880379973
- MSE test: 24.163437363863615

d. Comparison among different methods. (Average among 20 runs)

Method	MSE train	MSE test
Naive Regression	86.75±4.84	80.07±9.54
Linear Regression (attribute 1)	73.63±4.57	68.81±9.10
Linear Regression (attribute 2)	75.65±4.53	69.54±8.96
Linear Regression (attribute 3)	67.63±4.56	59.29±9.05
Linear Regression (attribute 4)	83.57±4.69	78.97±9.31
Linear Regression (attribute 5)	71.34±4.63	64.83±9.15
Linear Regression (attribute 6)	44.73±3.99	41.82±8.05
Linear Regression (attribute 7)	75.1±5.06	67.58±10.06
Linear Regression (attribute 8)	81.62±5.17	74.95±10.11
Linear Regression (attribute 9)	74.17±4.96	68.59±9.84
Linear Regression (attribute 10)	68.34±5.00	61.55±9.97
Linear Regression (attribute 11)	63.98±3.53	60.64±6.78
Linear Regression (attribute 12)	39.57±2.68	36.66±5.19
Linear Regression (all attributes)	22.49±2.07	23.77±4.68
Kernel Ridge Regression	8.44±1.65	12.15±2.14

2 Part II

To solve the questions in part II, we need to write a knn classifier first. Here we defined a knn classifier by sorting `euc_distance`, to improve its efficiency, we added `numba njit` (just in time compiler, no python mode) to enable fast compiling.

2.1 k -Nearest Neighbors

2.1.1 Generating the data

Question 6

In this question need to produce a figure of k -nn map trained from the sample data. The selected k for this example figure here is $k = 3$ and the test points is 100. To do this we need to complete the following tasks in this answer to question 6:

- Draw 100 random data uniformly from $[0, 1]^2$, name it X_Hs , which is the feature of 100 random points as the training set we draw; along with 100 random data uniformly from 0, 1, call it y_Hs , which is the label of 100 random points as the training set we draw. Those data would be inherited by answers to **question 7** and **question 8**, as the p_H distribution for following generating of train and test sets.
- Do k -NN predictions for the points on the set $[0, 1]^2$, in order to give label to each pixel on the background of the image plotted. This is done by meshgridding $[0, 1]^2$ by single pixel size of 0.001, which yields 1000000 pixels on the background to be labelled by the k nn-classifier with respect to the 100 training points generated in step 1 above.
- Plot the generated image, with 2 parts, firstly the background pixels to give a result of the k -NN classifier, next plot the 100 train points on the map. The result is shown below.

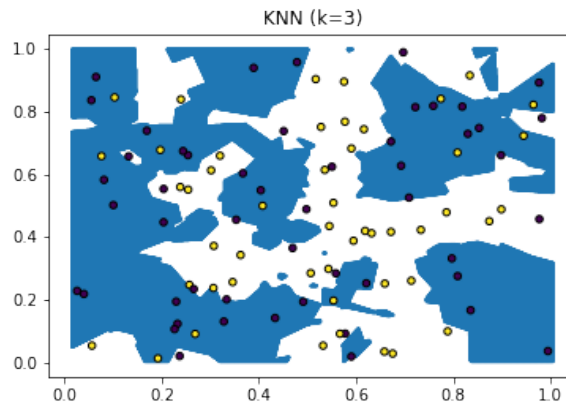


Figure 10: This figure is the visualization of generated hypothesis $h_{S,v}$ with $|S| = 100$ and $k = 3$. The white area is the mapping to 0 and the blue area is the mapping to 1, the corresponding centres are yellow and dark brown.

2.1.2 Estimated generalization error of k -NN as a function of k

Question 7

In this question there are 2 parts to be completed, firstly a figure would be given with respect to protocol A in the question listed. Secondly some comments would be given for this figure. This two parts could be combined and separated as the following tasks:

- Generate the random points following the instruction given in Question 7, as 80 percents of points would be drawn and given following the generated distribution p_H in question 6. While the other 20 percents points would be given by generating uniformly from 0, 1 as the noise in the data. This task is full-filled by writing a function called **"generating_points"** in the coding.
- Calculate the generalization error in the test sets. This is done by writing function called **"do_knn_test"**, each iteration of this function would produce a data containing the error calculated by comparing results of test labels and predicted labels. This test would be done for k from 1 to 49, each k would iterate 100 times and give the mean error of each 100 runs.
- Plot the figure. A figure with generalization error vs. k would be given. Comments would be attached in the description under the figure.

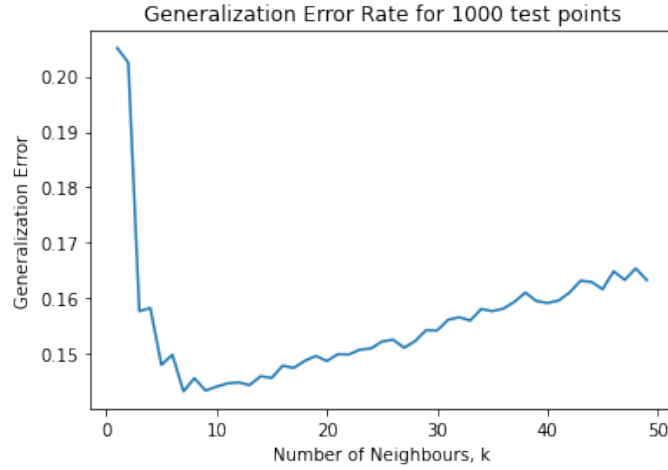


Figure 11: This figure is the visualization of generalization error vs. corresponding number of k .

- It could be seen that the value of generalization error was relatively high at $k = 1$, this is believed due to 1-NN Classifier would cause over-fitting, which caused the noise data in training set has been considered into the prediction. With k increasing to 10, the trend is sharply decreasing since the effect of over-fitting is decreasing.
- However, after k reached 10, the error would increase again with increment of k , that is because the majority of training and test points are generated from the distribution set in Question 6, which is 3-NN. So when k has exceeds too far from 3, larger k would lead to under-fitting, which means that too much details were lost in larger k -NN prediction

2.1.3 Determine the optimal k -NN as a function of the number of training points (m)

Question 8

This question could be separated as the following tasks, which is similar to question 7:

- Generate the random points following the instruction given in Question 7, as 80 percents of points would be drawn and given following the generated distribution p_H in question 6. While the other 20 percents points would be given by generating uniformly from 0, 1 as the noise in the data. This task is full-filled by writing a function called "**generating_points**" in the coding.
- Calculate the generalization error in the test sets. This is done by writing function called "**do_knn_test**", each iteration of this function would produce a data containing the error calculated by comparing results of test labels and predicted labels. This test would be done for k from 1 to 49, each k would iterate 100 times and give the mean error of each 100 runs. Minimum error would yield the selection of optimal k s for each iteration of m training points.
- Plot the figure. A figure with generalization error vs. k would be given. Comments would be attached in the description under the figure, same as done in question 7.

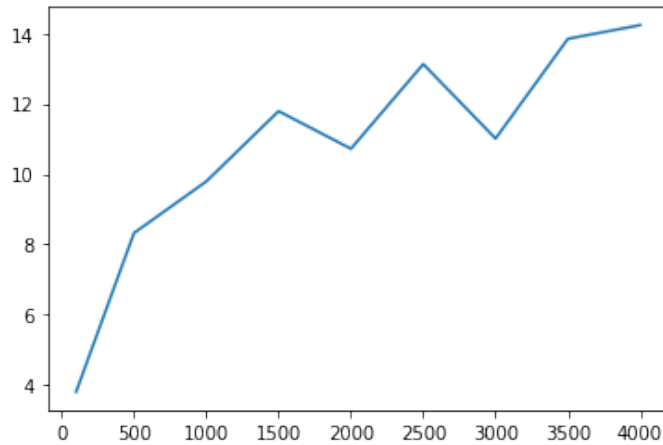


Figure 12: This figure is the visualization of optimal k vs. volume of training points. Here X-axis is the number of points in the training set. y-axis is the number of optimal k for each selection of training points.

- It could be seen that the value of optimal k is increasing with respect to the increment of numbers m of training points.
- The plotted line start from $[100, 3.8]$ which indicates that when the training set volume is the same as H_S, v generated in question 6, the optimal k is very close to 3.
- When the number of training points became larger, the optimal training points became larger since there are more noise points to be ignored in the set, so larger k would give robustness of k -NN model against increasing noise.

2.1.4 Additional Clarifications

This is done in Question 6 by adding a frame to the plot of question 6 in order to avoid corner points. By broaden the corner.

3 Part III

3.1 Questions

Question 9

Consider the function $K_c(\mathbf{x}, \mathbf{z}) := c + \sum_{i=1}^n x_i z_i$ where $\mathbf{x}, \mathbf{z} \in R^n$.

a) For what values of $c \in R$ is K_c a positive definite kernel?

The suitable value of c that would make K_c a positive definite kernel is $c \geq 0$.

Proof:

To prove that K_c is a PSD kernel, need to prove the gram matrix of K_c is a PSD Matrix.

The gram matrix of K_c is:

$$[K_c(x_i, x_j)]_{m \times m}$$

where $x_i, x_j \in \mathbf{x}$.

To prove the above matrix is a PSD matrix, refer to the property of PSD matrix, could rewrite the above matrix into the sum of two PSD matrix.

Define matrix $A =$

$$[< x_i, x_j >]_{m \times m}$$

where $x_i, x_j \in \mathbf{x}$, $< x_i, x_j >$ is the inner product of them. Here matrix A is a PSD matrix obviously.

Next define a Constant matrix $C =$

$$[c]_{m \times m}$$

where each term of matrix C is constant c in the kernel function K_c . To make the gram matrix of K_c a PSD, should ensure A and C PSD st.

$$K_c = A + B$$

Hence K_c a PSD Kernel if $c \geq 0$.

Which concludes the proof.

b) Suppose we use K_c as a kernel function with linear regression(least squares). Explain how c influences the solution.

Solve:

We define

$$K(\mathbf{x}, \mathbf{z}) := \sum_{i=1}^n x_i z_i, \text{ where } \mathbf{x}, \mathbf{z} \in R^n$$

$$K := \begin{bmatrix} K(\mathbf{x}_1, \mathbf{x}_1) & \dots & K(\mathbf{x}_1, \mathbf{x}_m) \\ \dots & \ddots & \dots \\ K(\mathbf{x}_m, \mathbf{x}_1) & \dots & K(\mathbf{x}_m, \mathbf{x}_m) \end{bmatrix}$$

Since,

$$K_c(\mathbf{x}, \mathbf{z}) := c + \sum_{i=1}^n x_i z_i, \text{ where } \mathbf{x}, \mathbf{z} \in R^n$$

Therefore,

$$K_c := \begin{bmatrix} K_c(\mathbf{x}_1, \mathbf{x}_1) & \dots & K_c(\mathbf{x}_1, \mathbf{x}_m) \\ \dots & \ddots & \dots \\ K_c(\mathbf{x}_m, \mathbf{x}_1) & \dots & K_c(\mathbf{x}_m, \mathbf{x}_m) \end{bmatrix} = C + K, \text{ where } C = \begin{bmatrix} c & \dots & c \\ \dots & \ddots & \dots \\ c & \dots & c \end{bmatrix}$$

when we apply $K(\mathbf{x}, \mathbf{z})$ into linear regression, we can calculate the weights

$$\alpha = K^{-1}y$$

Similarly, when we apply $K_c(\mathbf{x}, \mathbf{z})$ into linear regression, the weights is

$$\alpha_c = K_c^{-1}y = (C + K)^{-1}y$$

Hence,

$$\begin{aligned} K\alpha &= (C + K)\alpha_c \\ \alpha &= K^{-1}(C + K)\alpha_c = (K^{-1}C + I_m)\alpha_c \\ \alpha_c &= (K^{-1}C + I_m)^{-1}\alpha \end{aligned}$$

From the formulation above, we can see how c affects regression's solution α_c . The solution α_c can be taken over by c , If $c=0$, $C=O$, $\alpha_c = \alpha$.

Question 10

Suppose we perform linear regression with a Gaussian kernel $K_\beta(\mathbf{x}, \mathbf{t}) = \exp(-\beta\|\mathbf{x} - \mathbf{t}\|^2)$ to train a classifier on a dataset $(\mathbf{x}_1, y_1) \dots (\mathbf{x}_m, y_m) \in R^n \times \{-1, 1\}$. Thus obtaining a function $f : R^n \rightarrow R$ which is of the form $f(\mathbf{t}) = \sum_{i=1}^m \alpha_i K_\beta(\mathbf{x}_i, \mathbf{t})$. The corresponding classifier is then $\text{sign}(f(\mathbf{t}))$. This classifier depends on the parameter β selected for the kernel. In what scenario will chosen β enable the trained linear classifier to simulate a 1-NEAREST NEIGHBOUR CLASSIFIER trained on the same dataset?

Conclusion A relatively large β should be chosen to ensure that the kernel regression classifier could simulate the effect of 1-NN Classifier.

Proof Sketch a proof and list the chosen of β as a function of training feature \mathbf{x} and test feature \mathbf{t} . Firstly, the Gaussian Kernel could be rewritten as the form of amplifying a Gaussian Distribution pdf:

$$K_\beta(\mathbf{x}, \mathbf{t}) = \exp(-\beta\|\mathbf{x} - \mathbf{t}\|^2) = \left(\frac{1}{\sigma(\beta)\sqrt{2\pi}} \exp\left(\frac{-\|\mathbf{x} - \mathbf{t}\|^2}{2\sigma(\beta)^2}\right)\right) \times C(\beta)$$

where $C(\beta) = \sigma(\beta)\sqrt{2\pi}$, and $\sigma(\beta)$ defined implicitly by the following equation:

$$\beta = \frac{1}{2\sigma(\beta)^2}$$

It is obvious that

$$\frac{1}{\sigma(\beta)\sqrt{2\pi}} \exp\left(\frac{-\|\mathbf{x} - \mathbf{t}\|^2}{2\sigma(\beta)^2}\right)$$

forms a Gaussian distribution with mean at \mathbf{t} and with variance $\sigma(\beta)^2$.

From property of Gaussian distribution, more than 99% of possibilities laid in the interval $[\mathbf{t} - 3\sigma(\beta), \mathbf{t} + 3\sigma(\beta)]$, using this property, we define a $\sigma(\beta)$ that is small enough that allows only one point in \mathbf{x} laid in the interval $[\mathbf{t} - 3\sigma(\beta), \mathbf{t} + 3\sigma(\beta)]$. Which yields the following inequation:

$$[\min(\|\mathbf{x}_i - \mathbf{t}_i\|) | \text{first minimum element}] > 3\sigma(\beta) > [\min(\|\mathbf{x}_i - \mathbf{t}_i\|) | \text{second minimum element}]$$

where manually selected only 1 point x_i in \mathbf{x} to be covered in the range:

$$[\mathbf{t} - 3\sigma(\beta), \mathbf{t} + 3\sigma(\beta)]$$

Now, by selecting $\sigma(\beta)$ with respect to (\mathbf{x}, \mathbf{t}) , a relationship between β and (\mathbf{x}, \mathbf{t}) is constructed.

The affect of selecting β like above would yields only 1 points in $[\mathbf{x}]$ contribute more than 99% affect on the classifier function $f(\mathbf{t})$, which is just similar to the effect of 1-NN Classifier.

Which concludes the proof.

Question 11

note: all of the operations here observe the mod 2 operation rule. In the game of whack a mole, we have a $n \times n$ board, and moles come out from some holes in the board. We use a $n \times n$ matrix $B=[b_{ij}]$ to represent the board.

$$b_{ij} = \begin{cases} 1 & \text{mole comes out} \\ 0 & \text{mole hides} \end{cases}, i, j = 1 \dots n$$

Our purpose is to make $b_{ij} = 0$ for any $i, j = 1 \dots n$. For the question we know that every time we hit a mole with the mallet we'll cause the mole and the immediate four adjacent holes to change. So we can define the $\text{Hit}(p, q) = [H(p, q)_{ij}]$ as a $n \times n$ matrix, as below:

$$H(p, q)_{ij} = \begin{cases} 1 & (i, j) = (p, q) \text{ or } (p-1, q) \text{ or } (p+1, q) \text{ or } (p, q+1) \text{ or } (p, q-1) \\ 0 & \text{otherwise} \end{cases}, i, j = 1 \dots n$$

We give an initial board configuration as $B^{(0)}$ and we consider $B^{(k)}$ as the result of k-th hit. So every time we whack a mole in the position b_{pq} , we can get the iteration formulation below:

$$B^{(k+1)} = B^{(k)} + Hit(p, q)$$

If this problem exists a solution, we can find a finite sequence of $\{(p_m, q_k) | m = 1 \dots n; k = 1 \dots n\}$ to operate the $H(p_m, q_k)$ operation in order. So we can establish mathematical expressions, as below:

$$B_0 + \sum_{m=1}^n \sum_{k=1}^n \alpha_{m,k} H(p_m, q_k) = O, \alpha_{m,k} = 0 \text{ or } 1$$

$$\sum_{m=1}^n \sum_{k=1}^n \alpha_{m,k} H(p_m, q_k) = B_0$$

Therefore, the origin problem has been reduced to a problem of solving linear equations. Now we have n^2 equations (because there is a equation relationship for every elements in the matrix). The gaussian elimination method's complexity is $O(m^3)$, where m is the rank of linear equation system. In this problem, $m = n^2$, so the complexity of this problem based on our algorithm is $O(n^6)$.

A Appendix I: Code for Part I

A.1 Plot fit function curve, give corresponding equation and calculate mse for Q1

```

1 import numpy as np
2 from numpy.linalg import inv
3 import matplotlib.pyplot as plt
4 # define mse function
5 def mse(label_y, predict_y):
6     se_per = np.zeros_like(predict_y)
7     for index in range(len(label_y)):
8         se_per[index] = (label_y[index] - predict_y[index])**2
9     return np.mean(se_per)
10 #define regression class
11 class regression(object):
12     def __init__(self, input_x=[1,2,3,4], label_y=[3,2,0,5]):
13         self.input_x = input_x
14         self.label_y = label_y
15     #define linear regression function, return coefficients of basis functions, and MSE, and
    predict value
16
17     def polynomial_basis_fun_reg(self, order):
18         A = np.mat([0. for i in range(len(self.input_x)*order)]).reshape(len(self.input_x),
    order)
19         for row in range(len(self.input_x)):
20             for k in range(order):
21                 A[row, k] = (self.input_x[row])**k
22         b = np.mat(self.label_y).reshape(len(self.label_y), 1)
23         alpha = inv((A.T @ A) @ (A.T @ b)) #coefficiences of basis functions
24         #print("regression coefficients for order "+str(order)+" are:\n", alpha)
25         alpha = np.array(alpha).reshape(len(alpha))
26         #-----caulate MSE-----#
27         predict_y = np.array((A @ alpha))
28         predict_y = predict_y.reshape(len(self.label_y))
29         MSE = mse(self.label_y, predict_y)
30         #print("MSE for order "+str(order)+" is:', MSE)
31         return alpha, MSE, predict_y
32 #show the equations of fit function
33 def create_curves_equations(order, coefficients):
34     equation = ""
35     for index in range(order):
36         s = '{x^{}}'.format(coefficients[index, 0], index)
37         equation = equation + "+" + s
38     return equation[1:]
39 order = 4
40 create_curves_equations(order, coefficients4)

```

Listing 1: Code for Q1

A.2 Plot mse(or Ln(mse)) for both training and test set in Q2 to explore overfitting

```
1 #generate sample data
2 def generate_traning_x(sample_num=30,round=0):
3     np.random.seed(round)
4     input_x = np.random.uniform(0,1,sample_num)
5     return input_x
6 input_x=generate_traning_x()
7 #generate label data
8 def sin_square_sigma(x,sample_num=30,miu=0,sigma=0.07,round=0):
9     y=np.zeros_like(x)
10    np.random.seed(round)
11    epsilon=np.random.normal(miu,sigma,sample_num)
12    for i in range(len(x)):
13        y[i]=(np.sin(2*(np.pi)*x[i]))**2 + epsilon[i]
14    return y
15 #define a LN function to calculate ln value for every element in a list
16 def LN(error):
17     LN_result=np.zeros_like(error)
18     for i in range(len(error)):
19         LN_result[i]=np.log(error[i])
20     return LN_result
21 #(a)i plot scatter of our generating points and the fuction curve of (sin(2*pi*x))^2
22 x=np.linspace(0,1,100)
23 plt.plot(x, (np.sin(2*np.pi*x))**2)
24 plt.scatter(input_x,label_y,s=25,c='r')
25 plt.savefig('/SL_CW/pic/Q2ai.png')
26 #(a)ii
27 #define poly_k to calculate predict result of given x via our trained model
28 def poly_k(alpha,x,order):
29     y=0
30     for k in range(order):
31         y=y+alpha[k]*x**k
32     return y
33 #Fit the data set with a polynomial basis of dimension k=2,5,10,14,18 and plot each curve
34 def Q2aii():
35     input_x=generate_traning_x()
36     label_y=sin_square_sigma(input_x)
37     x=np.linspace(0,1,100)
38     w=[]
39     MSE=np.zeros((5,1))
40     #def plot_graph(input_x=[1,2,3,4],label_y=[3,2,0,5],max_order):
41     #predict_result=np.zeros((max_order,len(label_y)))#predict result for every order
42     test=regression(input_x,label_y)
43     alpha2,MSE[0],_=test.polynomial_basis_fun_reg(2)
44     alpha5,MSE[1],_=test.polynomial_basis_fun_reg(5)
45     alpha10,MSE[2],_=test.polynomial_basis_fun_reg(10)
46     alpha14,MSE[3],_=test.polynomial_basis_fun_reg(14)
47     alpha18,MSE[4],_=test.polynomial_basis_fun_reg(18)
48
49     #plot=np.poly1d(predict_result[0,:])
50     #
51     y2=np.zeros_like(x)
52     y5=np.zeros_like(x)
53     y10=np.zeros_like(x)
54     y14=np.zeros_like(x)
55     y18=np.zeros_like(x)
56     for i in range(len(x)):
57         y2[i]=poly_k(alpha2,x[i],2)
58         y5[i]=poly_k(alpha5,x[i],5)
59         y10[i]=poly_k(alpha10,x[i],10)
60         y14[i]=poly_k(alpha14,x[i],14)
61         y18[i]=poly_k(alpha18,x[i],18)
62     l2=plt.plot(x, y2,label="k=2")
63     l5=plt.plot(x, y5,label="k=5")
64     l10=plt.plot(x, y10,label="k=10")
65     l14=plt.plot(x, y14,label="k=14")
66     l18=plt.plot(x, y18,label="k=18")
67     plt.scatter(input_x,label_y,s=25,c='r')
68     # set up axis region
69     plt.legend(handles=[l2,l5,l10,l14,l18],labels=["k=2","k=5","k=10","k=14","k=18"])
70     plt.xlim((0, 1))
71     plt.ylim((0, 1))
```

```

72 plt.savefig('/SL_CW/pic/Q2a11.png')
73 plt.show()
74 #print(MSE)
75 Q2a11()
76 # (b)
77 #define class Error_analysis to analyse mse(or ln(mse)) both in training and test sets
78 class Error_analysis(object):
79     def __init__(self, input_x, label_y, max_order=18):
80         self.input_x=input_x
81         self.label_y=label_y
82         self.max_order=max_order
83         self.poly_dimension=np.arange(self.max_order+1)[1:]
84     #return training error
85     def calculate_mse(self):
86         #w=[]
87         training_error=np.zeros(self.max_order)
88         test_Q2=regression(self.input_x, self.label_y)
89         for k in range(self.max_order):
90             alpha, training_error[k], _=test_Q2.polynomial_basis_fun_reg(self.poly_dimension
91 [k])
92             w.append(alpha)
93         return training_error
94     #return ln(training error)
95     def calculate_Ln_mse(self):
96         training_error=np.zeros(self.max_order)
97         test_Q2=regression(self.input_x, self.label_y)
98         for k in range(self.max_order):
99             _, training_error[k], _=test_Q2.polynomial_basis_fun_reg(self.poly_dimension[k])
100         LN_training_error=LN(training_error)
101         return LN_training_error
102     #return Ln(test error) and test error
103     def ln_mse_in_testset(self, alpha, test_x, test_y, order):
104         A=np.mat([0. for i in range(len(test_x)*order)]).reshape(len(test_x), order)
105         for row in range(len(test_x)):
106             for k in range(order):
107                 A[row,k]=(test_x[row])**k
108             alpha=alpha.reshape(order,1)
109             predict_y=np.array(A@alpha)
110             predict_y=predict_y.reshape(len(test_y))
111             MSE=mse(test_y, predict_y)
112             ln_mse=np.log(MSE)
113             return ln_mse, MSE
114 #plot LN mse of training error against polynomial order k
115 def Q2b():
116     w=[]
117     input_x=generate_traning_x()
118     label_y=sin_square_sigma(input_x)
119     order_candidate=np.arange(19)[1:]
120     reg=regression(input_x, label_y)
121     for i in range(len(order_candidate)):
122         alpha, _, _=reg.polynomial_basis_fun_reg(order_candidate[i])
123         w.append(alpha)
124     Q2_b=Error_analysis(input_x, label_y)
125     LN_training_error_Q2_b=Q2_b.calculate_Ln_mse()
126     print(LN_training_error_Q2_b)
127     plt.xlim((0, 18))
128     plt.plot(order_candidate, LN_training_error_Q2_b)
129     plt.xlabel("k")
130     plt.ylabel("LN mse of training error")
131     plt.savefig('/SL_CW/pic/Q2b.png')
132     return w
133 w=Q2b()
134 # (c) plot LN mse of test error against polynomial order k
135 def Q2c(w):
136     input_x_Q2_c=generate_traning_x(sample_num=1000)
137     label_y_Q2_c=sin_square_sigma(input_x_Q2_c, sample_num=1000, miu=0, sigma=0.07)
138     Q2_c=Error_analysis(input_x_Q2_c, label_y_Q2_c)
139     LN_test_MSE=np.zeros((len(w),1))
140     order_candidate=np.arange(19)[1:]
141     for i in range(len(w)):
142         LN_test_MSE[i], _=Q2_c.ln_mse_in_testset(w[i], input_x_Q2_c, label_y_Q2_c,
143 order_candidate[i])
144     #Q2_c=Error_analysis(input_x_Q2_c, label_y_Q2_c)
145     #LN_training_error_Q2_c=Q2_c.calculate_Ln_mse()

```



```

220 def Q2d_plot():
221     order_candidate=np.arange(19)[1:]
222     line1,=plt.plot(order_candidate, LN_avg_te, label='train')
223     line2,=plt.plot(order_candidate, LN_avg_tse, label='test')
224     plt.legend(handles=[line1, line2], labels=["train", 'test'])
225     plt.xlim((0, 18))
226     plt.xlabel('order')
227     plt.ylabel('Ln average mse')
228     plt.savefig('/SL_CW/pic/Q2d.png')
229     plt.show()
230 Q2d_plot()

```

Listing 2: Q2 code

A.3 Plot mse(or $\ln(\text{mse})$) for both training and test set based on a different basis function for Q3

```

1 #define feature map based on sin(k*pi*x) basis function
2 def fm_sin(input_x, label_y, order):
3     m=len(input_x)
4     A=np.zeros((m,order))#
5     for row in range(m):
6         for col in range(order):
7             A[row,col]=np.sin((col+1)*np.pi*(input_x[row]))
8     return A
9 # (b)(c)
10 #plot ln average mse over only one round for training and test set
11 def Q3(r=1):
12     te=np.zeros((r,len(order_set)))
13     tse=np.zeros((r,len(order_set)))
14     for i in range(r):
15         b=i
16         if i==6:
17             b=107
18         elif i==21:
19             b=102
20         elif i==66:
21             b=108
22         elif i==95:
23             b=105
24         elif i==98:
25             b=104
26         train_x=generate_training_x(round=b)
27         train_y=sin_square_sigma(train_x, round=b)
28         test_x=generate_training_x(sample_num=1000, round=b)
29         test_y=sin_square_sigma(test_x, sample_num=1000, round=b)
30         w_all_orders=[]
31         #calculate training error
32         for k in range(len(order_set)):
33             A=fm_sin(train_x, train_y, order_set[k])
34             w=poly_reg(A, train_y)
35             w_all_orders.append(w)#reserve weights for test_set
36             predict_y=predict(A, w, order_set[k])
37             #print(predict_y)
38             te[i,k]=mse(predict_y, train_y)
39         for k in range(len(order_set)):
40             A_test=fm_sin(test_x, test_y, order_set[k])
41             predict_tsy=predict(A_test, w_all_orders[k], order_set[k])
42             #print(predict_tsy)
43             tse[i,k]=mse(predict_tsy, test_y)
44         LN_avg_te=LN(np.mean(te, axis=0))
45         LN_avg_tse=LN(np.mean(tse, axis=0))
46         return LN_avg_te, LN_avg_tse
47 LN_avg_te, LN_avg_tse=Q3()
48 #Do plot
49 def Q3_plot():
50     order_candidate=np.arange(19)[1:]
51     line1,=plt.plot(order_candidate, LN_avg_te, label='train')
52     line2,=plt.plot(order_candidate, LN_avg_tse, label='test')
53     plt.legend(handles=[line1, line2], labels=["train", 'test'])
54     plt.xlim((0, 18))
55     plt.xlabel('order')

```



```

56 plt.ylabel('Ln average mse')
57 plt.savefig('/SL_CW/pic/Q3bc.png')
58 plt.show()
59 Q3_plot()
60 #(d)plot ln average mse over 100 round for training and test set
61 def Q3(r=100):
62     te=np.zeros((r,len(order_set)))
63     tse=np.zeros((r,len(order_set)))
64     for i in range(r):
65         b=i
66         if i==6:
67             b=107
68         elif i==21:
69             b=102
70         elif i==66:
71             b=108
72         elif i==95:
73             b=105
74         elif i==98:
75             b=104
76
77     train_x=generate_traning_x(round=b)
78     train_y=sin_square_sigma(train_x,round=b)
79
80     test_x=generate_traning_x(sample_num=1000,round=b)
81     test_y=sin_square_sigma(test_x,sample_num=1000,round=b)
82
83     w_all_orders=[]
84     #calculate training error
85     for k in range(len(order_set)):
86         A=fm_sin(train_x,train_y,order_set[k])
87         w=poly_reg(A,train_y)
88         w_all_orders.append(w)#reserve weights for test_set
89         predict_y=predict(A,w,order_set[k])
90         #print(predict_y)
91         te[i,k]=mse(predict_y,train_y)
92     for k in range(len(order_set)):
93         A_test=fm_sin(test_x,test_y,order_set[k])
94         predict_tsy=predict(A_test,w_all_orders[k],order_set[k])
95         #print(predict_tsy)
96         tse[i,k]=mse(predict_tsy,test_y)
97     LN_avg_te=LN(np.mean(te,axis=0))
98     LN_avg_tse=LN(np.mean(tse,axis=0))
99     return LN_avg_te,LN_avg_tse
100 LN_avg_te,LN_avg_tse=Q3()
101 #Do plot
102 def Q3_plot():
103     order_candidate=np.arange(19)[1:]
104     line1=plt.plot(order_candidate,LN_avg_te,label='train')
105     line2=plt.plot(order_candidate,LN_avg_tse,label='test')
106     plt.legend(handles=[line1,line2],labels=["train","test"])
107     plt.xlim((0, 18))
108     plt.xlabel('order')
109     plt.ylabel('Ln average mse')
110     plt.savefig('/SL_CW/pic/Q3d.png')
111     plt.show()
112 Q3_plot()

```

Listing 3: Q3 code

A.4 Baseline versus full linear regression

```

1 #naive regression,linear_reg_single_attributes,linear_reg_all_attributes
2 class naive_regression(object):
3     def __init__(self,input_x,label_y):
4         self.input_x=input_x
5         self.label_y=label_y
6     def fm(self):
7         m=len(self.input_x)
8         A=np.ones((m,1))
9         return A
10    def naive_reg(self,A):
11        w=np.linalg.inv(A.T@A)@(A.T)@(self.label_y)
12        w=np.array(w)

```

```

13     w=w.reshape(len(w))
14     return w
15 def predict(self,A,w):
16     test_y=A@w
17     test_y=np.array(test_y)
18     test_y=test_y.reshape(len(test_y))
19     return test_y
20 class linear_reg_single_attributes(object):
21     def __init__(self,input_x,label_y,rank):
22         self.rank=rank
23         self.x=input_x[:,self.rank]
24         self.label_y=label_y
25     def fm(self):
26         m=len(self.x)
27         A=np.zeros((m,2))
28         for row in range(m):
29             for k in range(2):
30                 A[row,k]=(self.x[row])**k
31         return A
32     def single_reg(self,A):
33         w=np.linalg.inv(A.T@A)@(A.T)@(self.label_y)
34         w=np.array(w)
35         w=w.reshape(len(w))
36         return w
37     def predict(self,A,w):
38         test_y=A@w
39         test_y=np.array(test_y)
40         test_y=test_y.reshape(len(test_y))
41         return test_y
42 class linear_reg_all_attributes(object):
43     def __init__(self,input_x,label_y):
44         self.input_x=input_x
45         self.label_y=label_y
46     def fm(self):
47         m=len(self.input_x)
48         bias=np.ones((len(self.input_x),1))
49         A=np.hstack((bias,self.input_x))
50         return A
51     def multi_reg(self,A):
52         w=np.linalg.inv(A.T@A)@(A.T)@(self.label_y)
53         w=np.array(w)
54         w=w.reshape(len(w))
55         return w
56     def predict(self,A,w):
57         test_y=A@w
58         test_y=np.array(test_y)
59         test_y=test_y.reshape(len(test_y))
60         return test_y
61 #perform comparision
62 def Q4():
63     totalRuns=20
64     gamma_candidates=np.array([2**(i) for i in range(-40,-25)])
65     sigma_candidates=np.array([2**7,2**(7.5),2**(8),2**(8.5),2**(9),2**(9.5),2**(10),
66                                2**(10.5),2**(11),2**(11.5),2**(12),2**(12.5),2**(13)])
67     dataset=pd.read_csv(r'/content/Boston-filtered.csv').to_numpy()
68     m,n=dataset.shape
69     num_trainset=int(m*(2/3))
70     #for single reg
71     train_error_SR=np.zeros((totalRuns,12))
72     test_error_SR=np.zeros((totalRuns,12))
73     #for naive reg
74     train_error_NR=np.zeros((totalRuns,1))
75     test_error_NR=np.zeros((totalRuns,1))
76     train_error_ML=np.zeros((totalRuns,1))
77     test_error_ML=np.zeros((totalRuns,1))
78     #-----#
79     for i in range(totalRuns):
80         dataset2=np.random.permutation(dataset)
81         X=dataset2[:, :-1]
82         Y=dataset2[:, -1:]
83         train_X=X[0:num_trainset]
84         train_Y=Y[0:num_trainset]
85         test_X=X[num_trainset:]
86         test_Y=Y[num_trainset:]
87         #linear_reg with single attribute

```

```

87     for rank in range(12):
88         Q4_SR=linear_reg_single_attributes(train_X,train_Y,rank)
89         A_train=Q4_SR.fm()
90         w=Q4_SR.single_reg(A_train)
91         predict_y=Q4_SR.predict(A_train,w)
92         train_error_SR[i,rank]=mse(predict_y,train_Y)
93         #test
94         Q4_SR_test=linear_reg_single_attributes(test_X,test_Y,rank)
95         A_test=Q4_SR_test.fm()
96         predict_tsy=Q4_SR.predict(A_test,w)
97         test_error_SR[i,rank]=mse(predict_tsy,test_Y)
98     #naive reg
99     Q4_NR=naive_regression(train_X,train_Y)
100    A_train2=Q4_NR.fm()
101    w2=Q4_NR.naive_reg(A_train2)
102    predict_y2=Q4_NR.predict(A_train2,w2)
103    train_error_NR[i]=mse(predict_y2,train_Y)
104    #test
105    Q4_NR_test=naive_regression(test_X,test_Y)
106    A_test2=Q4_NR_test.fm()
107    predict_tsy2=Q4_NR.predict(A_test2,w2)
108    test_error_NR[i]=mse(predict_tsy2,test_Y)
109    #multi-linear-reg
110    Q4_ML=linear_reg_all_attributes(train_X,train_Y)
111    A_train3=Q4_ML.fm()
112    w3=Q4_ML.multi_reg(A_train3)
113    predict_y3=Q4_ML.predict(A_train3,w3)
114    train_error_ML[i]=mse(predict_y3,train_Y)
115    #test
116    Q4_ML_test=linear_reg_all_attributes(test_X,test_Y)
117    A_test3=Q4_ML_test.fm()
118    predict_tsy3=Q4_ML.predict(A_test3,w3)
119    test_error_ML[i]=mse(predict_tsy3,test_Y)
120
121    #linear_reg with single attribute:mean and std of results of 20 rounds
122    avg_trainError_SR=np.mean(train_error_SR,axis=0)
123    std_trainError_SR=np.std(train_error_SR,axis=0)
124    avg_testError_SR=np.mean(test_error_SR,axis=0)
125    std_testError_SR=np.std(test_error_SR,axis=0)
126    #Naive Regression:mean and std of results of 20 rounds
127    avg_trainError_NR=np.mean(train_error_NR)
128    std_trainError_NR=np.std(train_error_NR)
129    avg_testError_NR=np.mean(test_error_NR)
130    std_testError_NR=np.std(test_error_NR)
131    #multi-linear-reg:mean and std of results of 20 rounds
132    avg_trainError_ML=np.mean(train_error_ML)
133    std_trainError_ML=np.std(train_error_ML)
134    avg_testError_ML=np.mean(test_error_ML)
135    std_testError_ML=np.std(test_error_ML)
136    print('avg_trainError_NR:',avg_trainError_NR)
137    print('std_trainError_NR:',std_trainError_NR)
138    print('avg_testError_NR:',avg_testError_NR)
139    print('std_testError_NR:',std_testError_NR)
140    print('+-----+')
141    print('avg_trainError_SR:',avg_trainError_SR)
142    print('std_trainError_SR:',std_trainError_SR)
143    print('avg_testError_SR:',avg_testError_SR)
144    print('std_testError_SR:',std_testError_SR)
145    print('+-----+')
146    print('avg_trainError_ML:',avg_trainError_ML)
147    print('std_trainError_ML:',std_trainError_ML)
148    print('avg_testError_ML:',avg_testError_ML)
149    print('std_testError_ML:',std_testError_ML)
150    print('+-----+')
151    Q4()

```

Listing 4: Q4 code

A.5 Perform KRR ,find the best pair of gamma and sigma for KRR, compare the performance with other methods.

```

1 #derived data set

```

```

2 dataset=pd.read_csv(r'/SL_CW/Boston-filtered.csv').to_numpy()
3 #define gaussian kernel
4 def Gaussian_kernel(xi,xj,sigma):
5     n=np.linalg.norm((xi-xj),ord=2)
6     return np.exp(n**2/(-2*sigma**2))
7 #define gaussian kernel feature map
8 def fm_GK(input_x,sigma):
9     m=len(input_x)
10    A=np.zeros((m,m))
11    for row in range(m):
12        for col in range(m):
13            A[row,col]=Gaussian_kernel(input_x[row],input_x[col],sigma)
14    return A
15 #define predict value
16 def predict_GK(train_x,test_x,w,sigma):
17     m=len(train_x)
18     n=len(test_x)
19     A_test=np.zeros((m,n))
20     for row in range(m):
21         for col in range(n):
22             A_test[row,col]=Gaussian_kernel(train_x[row],test_x[col],sigma)
23     test_y=(A_test.T)@w
24     test_y=np.array(test_y)
25     test_y=test_y.reshape(len(test_y))
26     return test_y
27 #define five fold cross validation for training hyper parameters
28 def five_fold_cv(train_X,train_Y,gamma,sigma):
29     #sigma=2**7
30     val_set_size=int(0.2*len(train_X))
31     test_mse_each_fold=np.zeros((5,1))
32     for i in range(5):
33         if i==4:
34             V_X=train_X[i*val_set_size:]
35             V_Y=train_Y[i*val_set_size:]
36             T_X=train_X[:i*val_set_size]
37             T_Y=train_Y[:i*val_set_size]
38         else:
39             V_X=train_X[i*val_set_size:(i+1)*val_set_size]
40             V_Y=train_Y[i*val_set_size:(i+1)*val_set_size]
41             T_X=np.vstack((train_X[:i*val_set_size],train_X[(i+1)*val_set_size:]))
42             T_Y=np.vstack((train_Y[:i*val_set_size],train_Y[(i+1)*val_set_size:]))
43             #Q5a=kernel_methods(T_X,T_Y)
44             A=fm_GK(T_X,sigma)
45             w=kernel_reg(A,T_Y,gamma)
46             predict_y=predict_GK(T_X,V_X,w,sigma)
47             test_mse_each_fold[i]=mse(predict_y,V_Y)
48     avgtest_mse=np.mean(test_mse_each_fold)
49     return avgtest_mse
50 # (a)
51 #define a function to calculate average mse for every pair of hyperparameters and find the
    best pair
52 def avgmse_with_parameters(train_X,train_Y,gamma_candidates,sigma_candidates):
53     num_g=len(gamma_candidates)
54     num_s=len(sigma_candidates)
55     mse_for_each_par=np.zeros((num_g,num_s))
56     for g in range(num_g):
57         for s in range(num_s):
58             mse_for_each_par[g,s]=five_fold_cv(train_X,train_Y,gamma_candidates[g],
                sigma_candidates[s])
59     best_par=np.where(mse_for_each_par==np.min(mse_for_each_par))
60     best_gamma_index=best_par[0][0]
61     best_sigma_index=best_par[1][0]
62     return mse_for_each_par,best_gamma_index,best_sigma_index
63 gamma_candidates=np.array([2**(i) for i in range(-40,-25)])
64 sigma_candidates=np.array([2**7,2**(7.5),2**(8),2**(8.5),2**(9),2**(9.5),2**(10),2**(10.5),
    2**(11),2**(11.5),2**(12),2**(12.5),2**(13)])
65 mse_for_each_par,best_gamma_index,best_sigma_index=avgmse_with_parameters(train_X,train_Y,
    gamma_candidates,sigma_candidates)
66 print(mse_for_each_par)
67 print(best_gamma_index)
68 print(best_sigma_index)
69 # (b) plot the "cross-validation error" as a function of gamma and sigma
70 from mpl_toolkits.mplot3d import Axes3D
71 import matplotlib.pyplot as plt
72 import numpy as np

```

```

73 fig=plt.figure(figsize=(10.0,9.0))
74 #fig = plt.figure()
75 ax = fig.add_subplot(111, projection='3d')
76 xs=sigma_candidates
77 ys=gamma_candidates
78 X,Y=np.meshgrid(xs,ys)
79 zs=mse_for_each_par
80 ax.plot_surface(X,Y,zs)
81 ax.set_xlabel('sigma')
82 ax.set_ylabel('gamma')
83 ax.set_zlabel('mse')
84 plt.savefig('D:\SL_CW\pic\Q5b.png')
85 plt.show()
86 #(c)
87 #calculate the training and test mse based on the best pair of gamma and sigma
88 def Q4c():
89     gamma=gamma_candidates[best_gamma_index]
90     sigma=sigma_candidates[best_sigma_index]
91     np.random.seed(0)
92     dataset3 = np.random.permutation(dataset)
93     X=dataset3[:, :-1]
94     Y=dataset3[:, -1:]
95     train_X=X[0:num_trainset]
96     train_Y=Y[0:num_trainset]
97     test_X=X[num_trainset:]
98     test_Y=Y[num_trainset:]
99     A=fm_GK(train_X,sigma)
100    w=kernel_reg(A,train_Y,gamma)
101    predict_y=predict_GK(train_X,train_X,w,sigma)
102    train_mse=mse(predict_y,train_Y)
103    predict_tsy=predict_GK(train_X,test_X,w,sigma)
104    test_mse=mse(predict_tsy,test_Y)
105    print(train_mse)
106    print(test_mse)
107 Q4c()
108 #(d) performance comparing among baseline predicting, sinlge/multi linear regression and
    kernel method over 20 rounds
109 def Q5d():
110     totalRuns=20
111     gamma_candidates=np.array([2**(i) for i in range(-40,-25)])
112     sigma_candidates=np.array([2**7,2**(7.5),2**(8),2**(8.5),2**(9),2**(9.5),2**(10),
        2**(10.5),2**(11),2**(11.5),2**(12),2**(12.5),2**(13)])
113     dataset=pd.read_csv('r'/content/Boston-filtered.csv').to_numpy()
114     m,n=dataset.shape
115     num_trainset=int(m*(2/3))
116     #for single reg
117     train_error_SR=np.zeros((totalRuns,12))
118     test_error_SR=np.zeros((totalRuns,12))
119     #for naive reg
120     train_error_NR=np.zeros((totalRuns,1))
121     test_error_NR=np.zeros((totalRuns,1))
122     train_error_ML=np.zeros((totalRuns,1))
123     test_error_ML=np.zeros((totalRuns,1))
124     #-----#
125     #for kernel method
126     Store_best_par=np.zeros((totalRuns,2))
127     train_error_KM=np.zeros((totalRuns,1))
128     test_error_KM=np.zeros((totalRuns,1))
129     for i in range(totalRuns):
130         dataset2=np.random.permutation(dataset)
131         X=dataset2[:, :-1]
132         Y=dataset2[:, -1:]
133         train_X=X[0:num_trainset]
134         train_Y=Y[0:num_trainset]
135         test_X=X[num_trainset:]
136         test_Y=Y[num_trainset:]
137         #linear_reg with single attribute
138         for rank in range(12):
139             Q4_SR=linear_reg_single_attributes(train_X,train_Y,rank)
140             A_train=Q4_SR.fm()
141             w=Q4_SR.single_reg(A_train)
142             predict_y=Q4_SR.predict(A_train,w)
143             train_error_SR[i,rank]=mse(predict_y,train_Y)
144             #test
145             Q4_SR_test=linear_reg_single_attributes(test_X,test_Y,rank)

```

```

146     A_test=Q4_SR_test.fm()
147     predict_tsy=Q4_SR.predict(A_test,w)
148     test_error_SR[i,rank]=mse(predict_tsy,test_Y)
149     #naive reg
150     Q4_NR=naive_regression(train_X,train_Y)
151     A_train2=Q4_NR.fm()
152     w2=Q4_NR.naive_reg(A_train2)
153     predict_y2=Q4_NR.predict(A_train2,w2)
154     train_error_NR[i]=mse(predict_y2,train_Y)
155     #test
156     Q4_NR_test=naive_regression(test_X,test_Y)
157     A_test2=Q4_NR_test.fm()
158     predict_tsy2=Q4_NR.predict(A_test2,w2)
159     test_error_NR[i]=mse(predict_tsy2,test_Y)
160     #multi-linear-reg
161     Q4_ML=linear_reg_all_attributes(train_X,train_Y)
162     A_train3=Q4_ML.fm()
163     w3=Q4_ML.multi_reg(A_train3)
164     predict_y3=Q4_ML.predict(A_train3,w3)
165     train_error_ML[i]=mse(predict_y3,train_Y)
166     #test
167     Q4_ML_test=linear_reg_all_attributes(test_X,test_Y)
168     A_test3=Q4_ML_test.fm()
169     predict_tsy3=Q4_ML.predict(A_test3,w3)
170     test_error_ML[i]=mse(predict_tsy3,test_Y)
171     #kernel method
172     Store_best_par[i,0],Store_best_par[i,1]=find_best_par(train_X,train_Y,
gamma_candidates,sigma_candidates)
173     ga,si=Store_best_par[i,0],Store_best_par[i,1]
174     A_train4=fm_GK(train_X,si)
175     w4=kernel_reg(A_train4,train_Y,ga)
176     predict_y4=predict_GK(train_X,train_X,w4,si)
177     train_error_KM[i]=mse(predict_y4,train_Y)
178     #test
179     predict_tsy4=predict_GK(train_X,test_X,w4,si)
180     test_error_KM[i]=mse(predict_tsy4,test_Y)
181     #linear_reg with single attribute:mean and std of results of 20 rounds
182     avg_trainError_SR=np.mean(train_error_SR,axis=0)
183     std_trainError_SR=np.std(train_error_SR,axis=0)
184     avg_testError_SR=np.mean(test_error_SR,axis=0)
185     std_testError_SR=np.std(test_error_SR,axis=0)
186     #Naive Regression:mean and std of results of 20 rounds
187     avg_trainError_NR=np.mean(train_error_NR)
188     std_trainError_NR=np.std(train_error_NR)
189     avg_testError_NR=np.mean(test_error_NR)
190     std_testError_NR=np.std(test_error_NR)
191     #multi-linear-reg:mean and std of results of 20 rounds
192     avg_trainError_ML=np.mean(train_error_ML)
193     std_trainError_ML=np.std(train_error_ML)
194     avg_testError_ML=np.mean(test_error_ML)
195     std_testError_ML=np.std(test_error_ML)
196     #KM:mean and std of results of 20 rounds
197     avg_trainError_KM=np.mean(train_error_KM)
198     std_trainError_KM=np.std(train_error_KM)
199     avg_testError_KM=np.mean(test_error_KM)
200     std_testError_KM=np.std(test_error_KM)
201     print('avg_trainError_NR:',avg_trainError_NR)
202     print('std_trainError_NR:',std_trainError_NR)
203     print('avg_testError_NR:',avg_testError_NR)
204     print('std_testError_NR:',std_testError_NR)
205     print('+-+-----+-+')
206     print('avg_trainError_SR:',avg_trainError_SR)
207     print('std_trainError_SR:',std_trainError_SR)
208     print('avg_testError_SR:',avg_testError_SR)
209     print('std_testError_SR:',std_testError_SR)
210     print('+-+-----+-+')
211     print('avg_trainError_ML:',avg_trainError_ML)
212     print('std_trainError_ML:',std_trainError_ML)
213     print('avg_testError_ML:',avg_testError_ML)
214     print('std_testError_ML:',std_testError_ML)
215     print('+-+-----+-+')
216     print('avg_trainError_KM:',avg_trainError_KM)
217     print('std_trainError_KM:',std_trainError_KM)
218     print('avg_testError_KM:',avg_testError_KM)
219     print('std_testError_KM:',std_testError_KM)

```

```

220     print('best pairs of parameters:', Store_best_par)
221 Q5d()

```

Listing 5: Q5 code

B Appendix II: Code for Part II

B.1 K-NN Classifier and generating points

```

1  import numpy as np
2  from numba import njit, jit
3  import numba
4  from matplotlib import pyplot as plt
5
6  def Draw_Hs():
7      '''Draw a distribution for the following process firstly'''
8      n_points = 100
9
10     X_Hs = np.random.rand(n_points, 2)
11     y_Hs = np.random.choice([0, 1], n_points)
12
13     return X_Hs, y_Hs
14
15 @njit
16 def euc_dis(sample1, sample2):
17     '''
18     Euclidean distance between two sample points
19     sample1: A test sample. 2-tuple
20     sample2: A test sample. 2-tuple
21     '''
22     distance = np.sqrt(np.sum((sample1 - sample2)**2))
23     return distance
24
25 @njit
26 def get_distance(X, testInstance):
27     """
28     Use numba to accelerate the process of getting distances.
29     """
30     distances = [euc_dis(x, testInstance) for x in X]
31     return distances
32
33 def knn_classify(X, y, testInstance, k):
34     '''
35     Given a test data point testInstance, predict its label from knn classifier.
36     X: Data feature
37     y: Data label
38     testInstance: test sample
39     k: Number of neighbors chosen for one vote center.
40     '''
41     distances = get_distance(X, testInstance)
42     kneighbors = np.argsort(distances)[:k]
43     count = np.bincount(y[kneighbors])
44     predict_label = np.argmax(count)
45     return predict_label
46
47 def generate_points(X_Hs, y_Hs, sample_size, noise_size, k):
48     '''
49     '''
50     X_Hs_sample = np.random.rand(sample_size, 2)
51     y_Hs_sample = [knn_classify(X_Hs, y_Hs, data, k) for data in X_Hs_sample]
52
53     X_noise = np.random.rand(noise_size, 2)
54     y_noise = np.random.choice([0, 1], noise_size)
55
56     X = np.r_[X_Hs_sample, X_noise]
57     y = np.r_[y_Hs_sample, y_noise]
58
59     return X, y

```

Listing 6: K-nn Classifier

B.2 Do k-nn labelling and plot question 6 figure

```
1 X_Hs, y_Hs = Draw_Hs()
2 #Draw the distribution of H_S,v which would be used in q6, 7 and 8.
3 def question6(X_Hs, y_Hs):
4     # Uniformly draw some points as the training set.
5     n_points = 100
6
7     X = X_Hs
8     y = y_Hs
9
10    # Choosing K for the classifier
11    k=3
12
13    # Visualization
14    x_min, x_max = X[:, 0].min() - 0.01, X[:, 0].max() + 0.01 #Avoid corner case
15    y_min, y_max = X[:, 1].min() - 0.01, X[:, 1].max() + 0.01 #Avoid corner case
16    xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.001),
17                          np.arange(y_min, y_max, 0.001))
18
19    Z = [knn_classify(X, y, data, k) for data in np.c_[xx.ravel(), yy.ravel()]]
20
21    plt.scatter(xx,yy,Z)
22
23    plt.scatter(X[:, 0], X[:, 1], c=-y,
24               s=20, edgecolor='k')
25    plt.title('KNN (k=%d)'%k)
26
27    plt.show()
28
29 question6(X_Hs, y_Hs)
```

Listing 7: K-nn test

B.3 Do k-nn test and plot question 7 figure

```
1 def do_knn_test(X_hs, y_Hs, train_points, test_points,k):
2     # Setting volume of train and test sets
3     # Uniformly drawing the train set and test set
4     train_sample_size = int(0.8 * train_points)
5     train_noise_size = int(0.2 * train_points)
6
7     X_train, y_train = generate_points(X_Hs, y_Hs, train_sample_size, train_noise_size, 3)
8
9
10    test_sample_size = int(0.8 * test_points)
11    test_noise_size = int(0.2 * test_points)
12
13    X_test, y_test = generate_points(X_Hs, y_Hs, test_sample_size, test_noise_size, 3)
14
15    # Predictions
16    predictions = [knn_classify(X_train, y_train, data, k) for data in X_test]
17    # Check the accuracy of predictions.
18    errornum = np.count_nonzero((predictions==y_test)==False)
19    return errornum/y_test.shape[0]
20
21 k_upbd = 50 #Set the upper bound for k, actual value of k would reach k-1
22 neighbours=np.arange(k_upbd)
23 #Initiate generalized error
24 generalized_error=np.zeros(len(neighbours))
25 for k in neighbours[1:k_upbd]:
26     #Initiate error
27     error = np.zeros(100)
28     for i in np.arange(100):
29         #Do 100 runs for each k in neighbours list.
30         error[i] = do_knn_test(X_Hs, y_Hs, 4000, 1000, k)
31     generalized_error[k] = np.mean(error)
32     #plot(generalized_error vs. number of k)
33 plt.xlabel('Number of Neighbours, k')
34 plt.ylabel('Generalization Error')
35 plt.title('Generalization Error Rate for 1000 test points')
```



```
36 plt.plot(neighbours[1:k_upbd],generalized_error[1:k_upbd])
```

Listing 8: Calculating generalization error

B.4 Do k-nn test and plot question 8 figure

```
1 import numpy as np
2 from matplotlib import pyplot as plt
3 def question8(k_upbd):
4     #Set the upper bound for k, actual value of k would reach k_upbd-1
5     arrayaa = [100]
6     arraybb = np.arange(500, 4500, 500)
7     train_points = np.r_[arrayaa, arraybb] #array([ 100,  500, 1000, 1500, 2000, 2500,
8     3000, 3500, 4000])
9     neighbours=np.arange(k_upbd)
10    moptimal_k = np.zeros(np.size(train_points))
11    j = 0
12    for m in train_points:
13        ioptimal_k = np.zeros(100)
14        for i in np.arange(100):
15            # Initiate generalized error
16            generalized_error=np.zeros(len(neighbours)) # array[0..49]
17            generalized_error[0]=10000 # Avoid k=zero's error being the min value
18            for k in neighbours[1:k_upbd]:
19                generalized_error[k] = do_knn_test(X_Hs, y_Hs, m, 1000, k)
20                ioptimal_k[i] = np.where(generalized_error == np.min(generalized_error))[0][0]
21            # First optimal k for this one-time run.
22            moptimal_k[j] = np.mean(ioptimal_k)
23            j += 1
24    print(moptimal_k)
25    plt.plot(train_points, moptimal_k)
26    return moptimal_k
27 k_optimals49=question8(50)
```

Listing 9: Finding optimal k