

Russel Stuart Daries
UCABRSD

STUDENT

FIRST LETTER OF SURNAME
D

PERSONAL TUTOR

Coursework Coversheet

This document is the assessed coursework coversheet for all Computer Science modules conducted at UCL for this academic year. Please do the following when submitting coursework:

- Staple a completed and signed copy of this form to every piece of assessed coursework you submit for modules in the Department of Computer Science.
- avoid the use of document containers such as cardboard or plastic covers, document wallets, ring binders or folders (unless otherwise stated by the lecturer concerned).

We do not wish to discourage students from discussing their work with fellow students and collaborating in solving problems. However you should avoid allowing the collaborative phase to approach too close to a final solution which might make

it impossible for you to make your own distinctive intellectual contribution. The key point is that you must not present the results of another person's work "as though they were your own".
<http://www.ucl.ac.uk/current-students/study/plagiarism/>

UCL has now signed up to use a sophisticated detection system (JISC Turn-It-In) to scan work for evidence of plagiarism, and the Department intends to use this for assessed coursework. This system gives access to billions of sources worldwide, including websites and journals, as well as work previously submitted to the Department, UCL and other universities

The collection point for "returned marked" coursework will be in a public area. If you wish to have your coursework returned privately. You MUST collect your coursework as soon as you are informed. □

Submission Details

Please ensure that the details you give are accurate and completed to the best of your knowledge.

COURSEWORK 1

Dr. Mark Herbster

MODULE CODE

MODULE NAME

DEADLINE / 2017

Declaration

I have read and understood the UCL and Departmental statements and guidelines concerning plagiarism.
I declare that:

1. This submission is entirely my own unaided work.
2. Wherever published, unpublished, printed, electronic or other information sources have been used as a contribution or component of this work, these are explicitly, clearly and individually acknowledged by appropriate use of quotation marks, citations, references and statements in the text.
3. In preparing this coursework, I discussed the general approach to take with the person(s) named below, however the content of the submission is my own work alone:

Person(s) with whom I have discussed this coursework:

Nitish Mutha



For Departmental Office Use

DATE AND TIME RECEIVED

INITIALS OR STAMP

COMPGI01: Supervised Learning Coursework 1

Due on Monday, January 16, 2017

Dr. Mark Herbster

Russel Daries: UCABRSD
Nitish Mutha: UCABMUT

January 16, 2017

Contents

Exercise 1	7
(a)	7
(b)	7
(c)	7
(d)	8
Exercise 2	9
(a)	9
(b)	9
(c)	10
Exercise 3	11
(1)	11
(2)	12
Exercise 4	13
(a)	13
(b)	13
(c)	14
Exercise 5	15
(a)	15
(b)	15
(c)	16
(d)	16
Exercise 6	17
(a)	17
(b)	17
Exercise 7	18
Exercise 8	20
Exercise 9	21
(a)	21
(b)	21
(c)	23
Exercise 10	24
(a)	24
(b)	24
(c)	24
(1)	25
(2)	25
(d)	26
Question 1	27

Question 2	28
(a)	28
(b)	29
Question 3	33
(a)	33
(b)	33
Question 4	34
(a)	34
(b)	35
(c)	35
Appendices	37
Exercise 1: Code	38
ex1.m	38
Exercise 2: Code	41
ex2.m	41
Exercise 4: Code	44
ex4.m	44
Exercise 5: Code	48
ex5.m	48
Exercise 6: Code	55
ex6.m	55
Exercise 7: Code	59
ex7.m	59
ex7.4.m	62
ex7.5.m	63
ex7.6.m	64
Exercise 9: Code	65
Ex9_a.m	65
Ex9_b.m	66
Ex9_c.m	68
Exercise 10: Code	70
Ex10_RR.m	70
dualcost.m	72
kriddereg.m	72
Question 4 (Part II): Code	73
question4_1.m (least squares)	73
question4_2.m (perceptron)	74
question4_3.m (winnow)	75
question4_4.m (1NN)	77
Library Functions	79
calculateLSR.m	79
calculateLSRex2.m	79
calculateLSRex4.m	80
calculateLSRex5.m	80
calculateLSRex6.m	81
calculateLSRTTest.m	81
meanSquareError.m	82

standard_deviation.m	82
calK.m	82
learnModel.m	83

List of Figures

1	MSE on training data.	8
2	MSE on test data.	8
3	Weight vector \mathbf{w} for 200 trials	8
4	MSE for train data ($\mathbf{x} \in \mathbb{R}^{10}$).	10
5	MSE for test data ($\mathbf{x} \in \mathbb{R}^{10}$).	10
6	Ridge Regression with 100 training samples.	13
7	Ridge Regression with 10 training samples.	14
8	Ridge Regression with average errors for 100 training samples.	14
9	Ridge Regression with average errors for 10 training samples.	14
10	MSEs for training(80), validation(20) and test data.	15
11	MSEs for training(8), validation(2) and test data.	15
12	MSEs for optimal gamma selection for 200 trials.	16
13	Ridge Regression with average errors for 100 training samples.	16
14	Ridge Regression with average errors for 10 training samples.	16
15	Cross validation with 100 training samples.	17
16	Cross validation with 10 training samples.	17
17	Optimal γ for 100 training samples with tuning methods.	18
18	Optimal γ for 10 training samples with tuning methods.	18
19	MSE on test data.	19
20	MSE on test data.	19
21	Average MSE for 13 attributes.	21
22	Box plot for MSE on the training set.	21
23	Box plot for MSE on the test set.	21
24	Box plot for MSE on the training set with all attributes.	23
25	Box plot for MSE on the test set with all attributes.	23
26	Cross validation error with Kernel Ridge Regression.	25
27	A separating hyperplane (\mathbf{w}, b) for a two dimensional set [1].	29
28	Estimated sample complexity of <i>Perceptron</i> algorithm.	34
29	Estimated sample complexity of <i>Winnow</i> algorithm.	34
30	Estimated sample complexity of <i>Least Squares</i> algorithm.	34
31	Estimated sample complexity of <i>1-Nearest Neighbour</i> algorithm.	34

List of Tables

1	LSR errors for \mathbf{w}^* with 100 training samples.	7
2	LSR errors for \mathbf{w}^* with 10 training samples.	7
3	Average Mean Squared errors for varying training sample size repeated 200 times.	8
4	Average \mathbf{w}^* for 200 trials.	8
5	LSR errors for \mathbf{w}^* with 100 training samples.	9
6	LSR errors for \mathbf{w}^* with 10 training samples.	9
7	Average Mean Squared errors for varying training sample size repeated 200 times.	9
8	Comparison of γ tuning methods.	18
9	Comparison of average γ	19
10	Comparison of MSE performances with naive regression.	21
11	Comparison of attribute MSE performances.	22
12	Comparison of attribute MSE performances with all attributes.	23
13	MSE for optimal parameters γ^* and σ^*	25
14	Comparison of regression methods.	26
15	Computational time for algorithms.	35

Exercise 1

(a)

The Least Squares Regression was applied to Exercise 1 in order to minimize the mean of the squared errors on all training samples. This can be expressed as finding the optimal weight vector \mathbf{w}^* as shown in the Equation 1 below.

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \frac{1}{l} \sum_{i=1}^l (\mathbf{x}'_i \mathbf{w} - y_i)^2 \quad (1)$$

The relevant $\mathbf{X}, \mathbf{y}, \mathbf{n}$ and \mathbf{w} vectors were created as required for Exercise 1 as shown in below.

```

1 %% Exercise 1(a)
2 w = randn(1,1);
3 x = randn(600,1);
4 n = randn(600,1);
5 y = x*w + n;

```

(b)

After all the required data matrices and vectors were created, using the Normal Equation as shown in Equation 2 below, the optimal \mathbf{w} could be calculated.

$$\mathbf{w}^* = (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{y} \quad (2)$$

The optimal \mathbf{w}^* and mean squared errors on the train and test sets can be seen in Table 1 below.

\mathbf{w}^*	Train Error	Test Error
0.5413	0.9601	1.0475

Table 1: LSR errors for \mathbf{w}^* with 100 training samples.

It can be seen from Table 1 above that the test error was larger than the training error, with a 100 training examples and 500 test examples. The weight vector \mathbf{w} did not generalise well.

(c)

The same experiment as discussed in Exercise 1(b) was repeated with only 10 training samples and 500 test samples taken from standard normal distribution. The resultant errors and \mathbf{w}^* can be seen in Table 2 below.

\mathbf{w}^*	Train Error	Test Error
0.7475	0.7584	1.0700

Table 2: LSR errors for \mathbf{w}^* with 10 training samples.

It can be seen that the training error in Table 2 was lower, but this can be attributed to only having sampled 10 times. The test error was higher by 0.0225 and the \mathbf{w}^* was larger by 0.2061.

(d)

The process followed in parts (a),(b) and (c) were repeated 200 times in order to evaluate the effect of varying the training size. Table 3 shown below demonstrated the average mean square errors for the training and test sets. Even though the training set performed well the test set performed poorly.

No. Training Samples	Train Error	Test Error
10	0.9247	1.1218
100	0.9921	1.0061

Table 3: Average Mean Squared errors for varying training sample size repeated 200 times.

Figure 1 and 2 shown below demonstrated the resultant MSE for each of the 200 trials on the training and the test set for the two training set sizes.

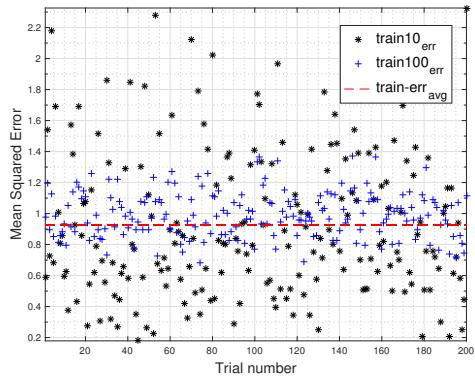


Figure 1: MSE on training data.

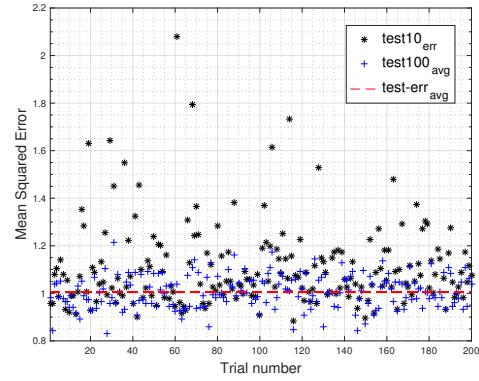


Figure 2: MSE on test data.

The average \bar{w}^* over the 200 experiments can be seen in Table 4 below.

10 Training Samples	100 Training Samples
0.0057	-0.0222

Table 4: Average w^* for 200 trials.

Figure 3 below demonstrated the calculated w^* for each trial.

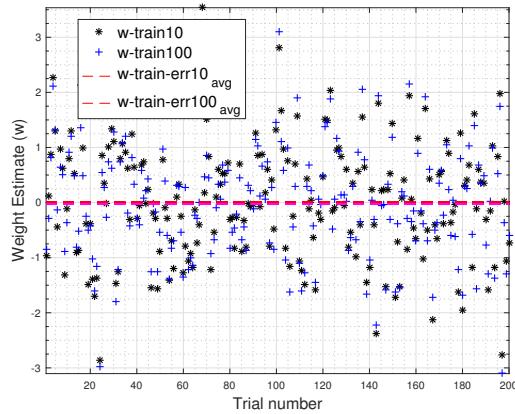


Figure 3: Weight vector w for 200 trials

Exercise 2

(a)

It could be seen in Exercise 1, the effect of the training and test corpus sizes was investigated in order to see how they affected the performance of the predictor. We will now extend the dimension of the data to $\mathbf{x} \in \mathbb{R}^{10}$.

The relevant $\mathbf{X}, \mathbf{y}, \mathbf{n}$ and \mathbf{w} vectors were created as required for Exercise 2(a) as shown in below.

```

1 %% Exercise 1A
2 w = randn(10,1);
3 x = randn(600,10);
4 n = randn(600,1);
5 y = x*w + n;

```

The training and test sets were created and tested using the *calculateLSRex2* function which can be found in Appendix ??.

(b)

A similar process as discussed in Exercise 1 (b)-(d) was followed now with a higher dimensional \mathbf{x} . The results for the experiments can be seen in Tables 5 - 7 shown below.

Train Error	Test Error
0.9367	1.1269

Table 5: LSR errors for \mathbf{w}^* with 100 training samples.

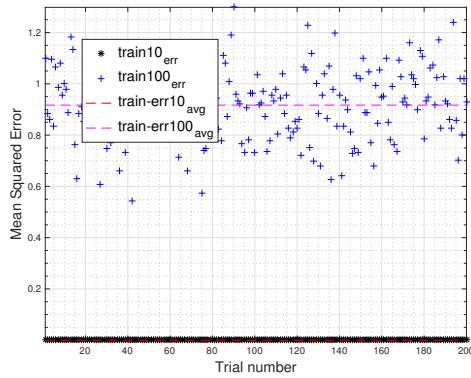
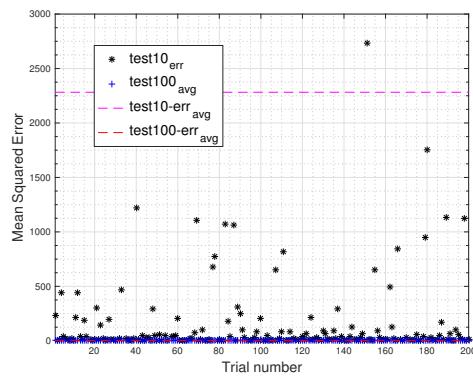
Train Error	Test Error
0.00	129.0357

Table 6: LSR errors for \mathbf{w}^* with 10 training samples.

Training Samples	Train Error	Test Error
10	0.0	2282.20
100	0.9170	1.1053

Table 7: Average Mean Squared errors for varying training sample size repeated 200 times.

The MSEs for each experiment can be seen in Figure 4 and 5 on the following page.

Figure 4: MSE for train data ($\mathbf{x} \in \mathbb{R}^{10}$).Figure 5: MSE for test data ($\mathbf{x} \in \mathbb{R}^{10}$).

The code utilized to perform these calculations can be seen in the Appendix.

(c)

It can be seen when comparing Figures 4 and 5 that the problem of overfitting is present. In the case of 10 training examples, the MSE was small 0.00 in comparison to MSE with 100 training examples with MSE 0.9170 . The resultant weight vector \mathbf{w}_{10} was overfitted, so when it was tested on the 500 test examples the resultant average MSE was 2282.20 as compared to the \mathbf{w}_{100} which generalised far better as the MSE of 1.1053 was much lower.

Thus is can be seen, the effect of the training corpus size of 10 examples is suseptible to fitting noise contained in the signal of interest resulting in low training errors but large test errors. did not result in a large test error. This is clearly evident when comparing the test MSE's shown in Table 7. Whereas, with a substantial amount of training examples, which in this case was 100, \mathbf{w}^* generalised far better when compared to 10 training examples even at higher dimensions of \mathbf{x} . By increasing the training examples to 100, however did not lead to a good indication of the performance of the regression function on the test set.

This is as a result of the regression function fitting the noise present in the signal of interest (overfitting).

Exercise 3

(1)

The optimization problem for Regularized Regression was solved in the following manner. Using Equation 1, Equation 3 could be proved.

$$\mathbf{w}^* = (\mathbf{X}'\mathbf{X} + \gamma l\mathbf{I})^{-1}\mathbf{X}'\mathbf{y} \quad (3)$$

The cost function $R(\mathbf{w})$ can be defined using Equation 1.

$$R(\mathbf{w}) = \frac{1}{l} \sum_{i=1}^l (\mathbf{x}'_i \mathbf{w} - y_i)^2 + \gamma \mathbf{w}' \mathbf{w} \quad (4)$$

Equation 4 is minimised when $\nabla_{\mathbf{w}} R(\mathbf{w}) = 0$. Therefore, the following derivation was performed.

$$\begin{aligned} \nabla_{\mathbf{w}} R(\mathbf{w}) &= \nabla_{\mathbf{w}} \left(\frac{1}{l} \sum_{i=1}^l (\mathbf{x}'_i \mathbf{w} - y_i)^2 + \gamma \mathbf{w}' \mathbf{w} \right) \\ &= \frac{1}{l} \sum_{i=1}^l \nabla_{\mathbf{w}} ((\mathbf{x}'_i \mathbf{w} - y_i)^2) + \nabla_{\mathbf{w}} (\gamma \mathbf{w}' \mathbf{w}) \\ &= \frac{1}{l} \nabla_{\mathbf{w}} (\mathbf{w}' \mathbf{X}' \mathbf{X} \mathbf{w} - 2\mathbf{y}' \mathbf{X} \mathbf{w} + \mathbf{y}' \mathbf{y}) + \nabla_{\mathbf{w}} (\gamma \mathbf{w}' \mathbf{w}) \\ \nabla_{\mathbf{w}} R(\mathbf{w}) &= \frac{1}{l} (2\mathbf{X}' \mathbf{X} \mathbf{w} - 2\mathbf{X}' \mathbf{y}) + (2\gamma \mathbf{w}) \end{aligned}$$

Then we set the derivative of the cost function to zero as shown below and solve the equation to make \mathbf{w} the subject of the formula.

$$\begin{aligned} \nabla_{\mathbf{w}} R(\mathbf{w}^*) &= 0 \\ \frac{1}{l} (2\mathbf{X}' \mathbf{X} \mathbf{w}^* - 2\mathbf{X}' \mathbf{y}) + (2\gamma \mathbf{w}^*) &= 0 \\ (\mathbf{X}' \mathbf{X} \mathbf{w}^* - \mathbf{X}' \mathbf{y}) + (\gamma l \mathbf{w}^*) &= 0 \\ \mathbf{X}' \mathbf{X} \mathbf{w}^* + \gamma l \mathbf{w}^* &= \mathbf{X}' \mathbf{y} \\ (\mathbf{X}' \mathbf{X} + \gamma l \mathbf{I}) \mathbf{w}^* &= \mathbf{X}' \mathbf{y} \\ \mathbf{w}^* &= (\mathbf{X}' \mathbf{X} + \gamma l \mathbf{I})^{-1} \mathbf{X}' \mathbf{y} \end{aligned}$$

Therefore, as shown in the working above the optimization problem yields the following equation.

$$\mathbf{w}^* = (\mathbf{X}' \mathbf{X} + \gamma l \mathbf{I})^{-1} \mathbf{X}' \mathbf{y} \quad (5)$$

(2)

The proof demonstrating the fact that the $\mathbf{X}'\mathbf{X} + \gamma l\mathbf{I}$ is positive definite matrix can be seen below.

The expression $\mathbf{X}'\mathbf{X} + \gamma l\mathbf{I}$ can be broken into two separate terms. For simplicity we will deal with the second term now.

$$\gamma l\mathbf{I}$$

Let matrix \mathbf{S} be defined as positive semi-definite if $\mathbf{x}^T \mathbf{S} \mathbf{x} \geq 0$ for all vectors $\mathbf{x} \neq 0$. Therefore, using this definition and the fact that \mathbf{S} is an identity matrix $\mathbf{S} = \mathbf{I}$, it is easy to prove that \mathbf{I} is positive definite for a 3×3 matrix that can be extended to a $n \times n$ matrix.

$$\mathbf{x}^T \mathbf{S} \mathbf{x} = [a \ b \ c] \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} a \\ b \\ c \end{bmatrix} = a^2 + b^2 + c^2 > 0 \quad (6)$$

Assuming the following inequalities hold true namely $\gamma > 0$ and $l > 0$. It is possible to state that the second term $\gamma l\mathbf{I}$ is positive definite.

The first term $\mathbf{X}'\mathbf{X}$ can be proved to be positive semi-definite as follows.

$$\mathbf{z}^T \mathbf{X}^T \mathbf{X} \mathbf{z} > 0$$

$$\mathbf{X} \mathbf{z}^T \mathbf{X} \mathbf{z} > 0$$

$$\|\mathbf{X} \mathbf{z}\|_2^2 > 0$$

Therefore, combining the two results obtained above the following statement can be made:

$$\mathbf{X}'\mathbf{X} + \gamma l\mathbf{I} > 0$$

Exercise 4

(a)

The weakness of normal regression was highlighted through the results presented in Exercises 1 and 2. However, through regularisation the freedom of the classifier is restricted such that the regression function will be less susceptible to noise and able to generalise better.

This is achieved by including a complexity term in the cost function. This is an adaptation of Equation 1 mentioned earlier. The new equation describing the optimal weight vector \mathbf{w}^* can be seen in Equation 7 shown below.

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \frac{1}{l} \sum_{i=1}^l (\mathbf{x}'_i \mathbf{w} - y_i)^2 + \gamma \mathbf{w}' \mathbf{w} \quad (7)$$

The regularisation parameter γ was varied over a range of values in order to investigate its effect on the performance of the regression function. Figure 6 shown below demonstrate training and test set errors as a function of γ .

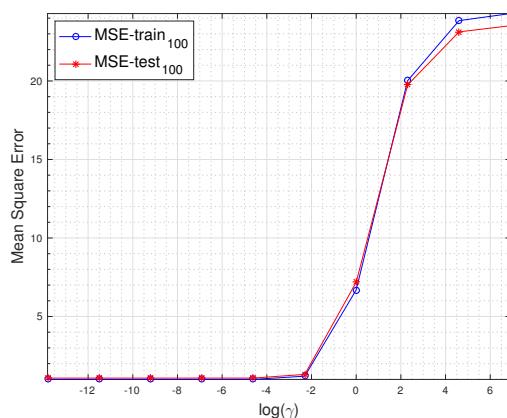


Figure 6: Ridge Regression with 100 training samples.

It can be seen from Figure 6 that the MSE of the test set is slightly higher than that of the training set but the regression function still generalises well with only 100 training samples.

(b)

The same process was followed but with only 10 training samples and the resultant MSE's for the training and test sets can be seen in Figure 7 shown on the following page.

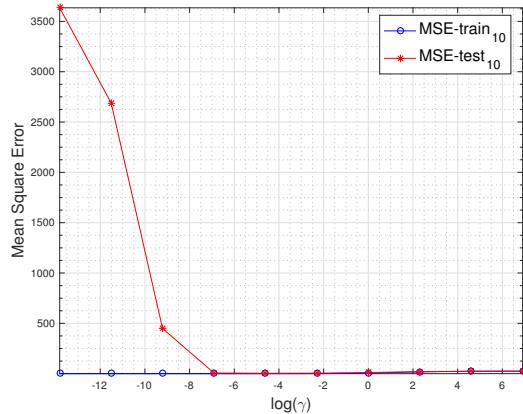


Figure 7: Ridge Regression with 10 training samples.

It can be seen that the phenomena of overfitting is still prevalent as the MSE for the training set is much lower when compared to the test set MSE.

(c)

Following this result, the process discussed in Exercise 4(a) and 4(b) was repeated for 200 independent trials. The average of the MSEs for each γ value for 200 trials was calculated in order to evaluate the performance. Figures 8 and 9 shown below demonstrates the results found.

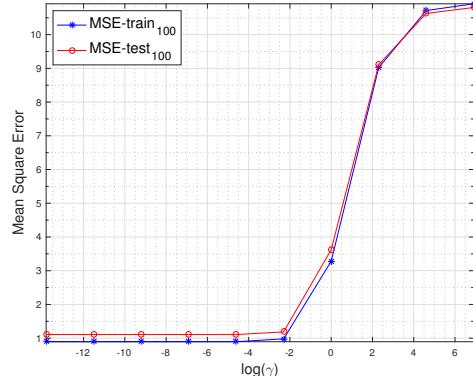


Figure 8: Ridge Regression with average errors for 100 training samples.

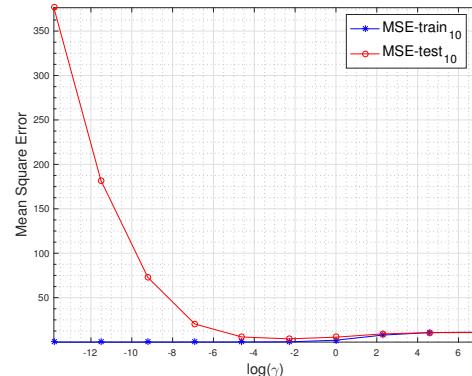


Figure 9: Ridge Regression with average errors for 10 training samples.

It can be seen from Figure 8 and 9 the general behaviour of the ridge regression over 200 trials remains matches the pattern identified in Exercise 4(a) and 4(b). However, just identifying the training set error and selecting the γ that minimizes the MSE is not sufficient in predicting the performance of the regression function on the test set. It is a poor indicator of the expected performance.

Exercise 5

(a)

As discussed in the previous exercise, only using the training set to select the best γ was insufficient. Therefore, as a result we will now use the validation set to select the best γ value and perform training on the complete training set (100 samples) which include the validation set (20 samples) for 200 trials and average the results.

The regularization value γ was varied over the validation set as shown in Figure 10 shown below.

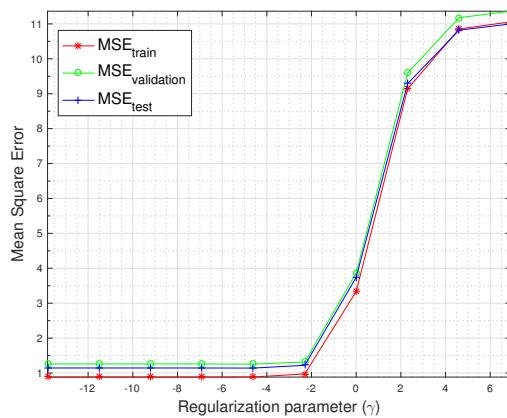


Figure 10: MSEs for training(80), validation(20) and test data.

The optimal regularizaton parameter $\bar{\gamma}^{(100)}$ was found to be 0.0445.

(b)

A similiar process for selection of γ was utilized as mentioned in Exercise 5(a) but with 8 training samples and 2 validation samples for 200 trials. The resultant plot can be seen in Figure 11 below.

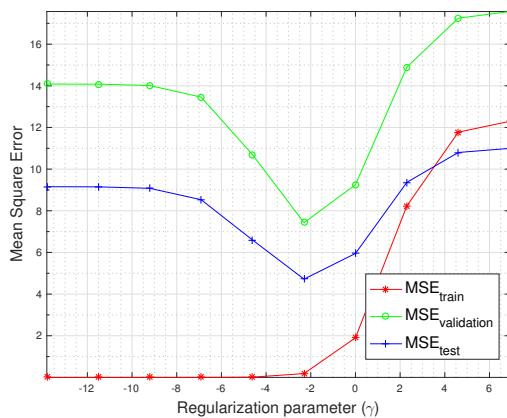


Figure 11: MSEs for training(8), validation(2) and test data.

The optimal regularizaton parameter $\bar{\gamma}^{(10)}$ was found to be 77.0543.

(c)

It can be seen from Figure 12 below that the MSE on the test set for 10 training samples even with the optimal gamma value is larger than that of 100 training samples.

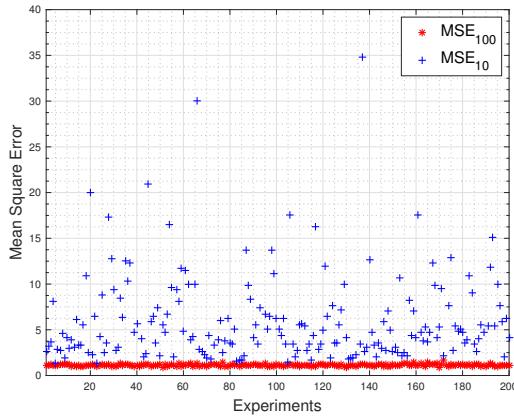


Figure 12: MSEs for optimal gamma selection for 200 trials.

It can also be seen that $\bar{\gamma}^{(10)} > \bar{\gamma}^{(100)}$ is larger because the complexity has larger weighting (Equation 7) in order to penalise the possibility of overfitting is greater for the 10 training samples.

(d)

A similar trend as that observed in Exercise 4(a) and 4(b) was observed when the dimensions of \mathbf{x} was reduced to one dimension. The results demonstrating this trend can be seen below in Figure 13 and 14.

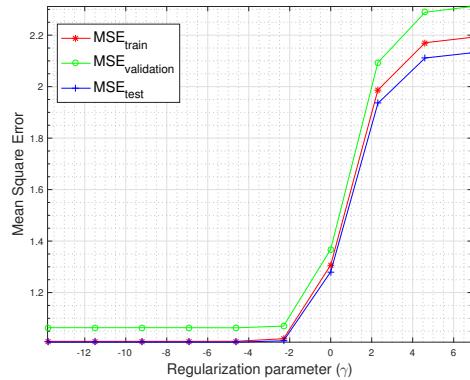


Figure 13: Ridge Regression with average errors for 100 training samples.

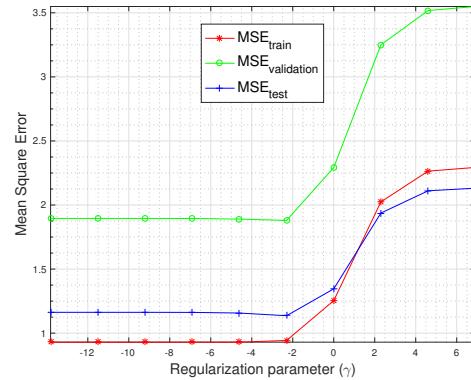


Figure 14: Ridge Regression with average errors for 10 training samples.

The validation MSE was considerably higher as can be seen in Figure 14, where as the test error was considerably less. The curves shown in Figure 13 matched relatively closely.

Exercise 6

(a)

The technique of using 5-fold cross-validation was used in order to select the optimal regularization parameter γ . As shown below the resultant cross-validation scores can be seen in Figure 15 with a training set size of 100.

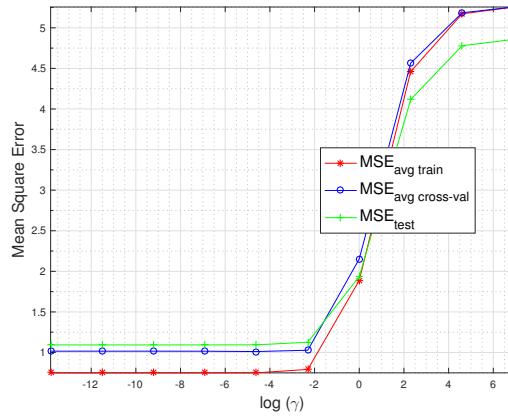


Figure 15: Cross validation with 100 training samples.

It can be seen from Figure 15 above the average cross-validation error score is slightly higher than the actual test set, suggesting the cross-validation method given enough samples, is able to generalise well. This is a similar behaviour to other selection methods presented in Exercise 5 and 6.

(b)

A similar process was followed for a training set size of 10. The result can be seen below in Figure 16.

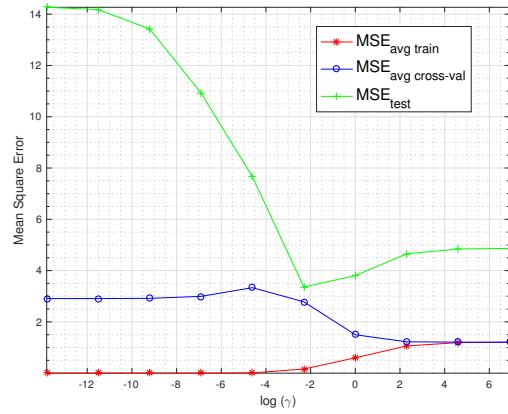


Figure 16: Cross validation with 10 training samples.

Analysing Figure 16 above, we can see that the test set error is larger than both the training and cross-validation scores. The trained regression function is a poor predictor. This can be attributed to the regression approach having insufficient data (10 samples) in order to perform well.

Exercise 7

Table 8 shown below demonstrates the comparison of γ tuning methods as covered in Exercises 4-6.

Tuning Method	No. of Training Samples	Average Test Error	Standard Deviation
Min Train Err.	10	2061.70	15938.00
Min Valid Err.	10	243.3223	2382.10
Cross Validate	10	6.1998	9.4548
Min Train Err.	100	1.1070	0.0864
Min Valid Err.	100	1.1200	0.0918
Cross Validate	100	1.1057	0.0854

Table 8: Comparison of γ tuning methods.

It is clearly evident from Table 8 that regardless of the optimal choice of γ the MSE on the test set given 10 training samples remains high with a large standard deviation.

The MSE on the test set with 100 training samples in Table 8 is far less with a very low standard deviation approximately 0.085 across all 3 tuning methods. The cross-validation method did achieve the best result on the test set followed by minimizing the training error method. The cross-validation method did achieve the best results due to the data segmentation allowing it to generalise better. The minimization of the training error makes the weights susceptible to over-fitting which is undesirable.

The mean error values for the optimal gamma selected for each of the tuning methods can be seen in Figure 17 and 18 can be seen below. It can be seen from these figures that the cross-validation error usually falls between the bounds of the training and validation errors.

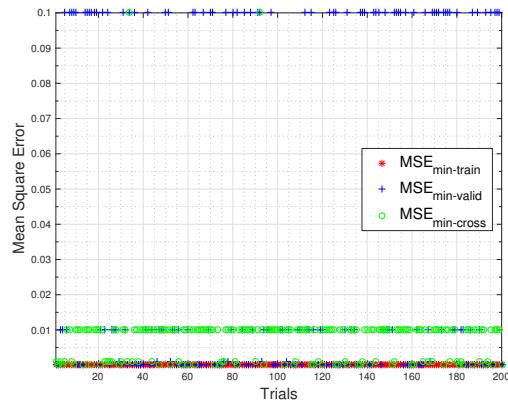


Figure 17: Optimal γ for 100 training samples with tuning methods.

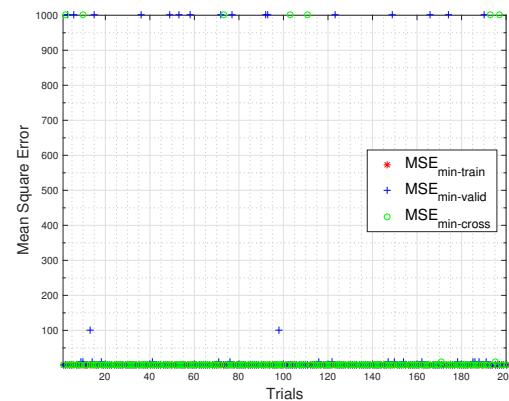


Figure 18: Optimal γ for 10 training samples with tuning methods.

Table 9 below shows the average γ over 200 trials.

Tuning Method	No. of Training Samples	Average γ
Min Train Err.	10	1×10^{-6}
Min Valid Err.	10	82.11
Cross Validate	10	35.28
Min Train Err.	100	1×10^{-6}
Min Valid Err.	100	0.0365
Cross Validate	100	0.0076

Table 9: Comparison of average γ .

The relevant scatter plots for each of the three tuning methods with training samples can be seen in Figure 19 (100 samples) and 20 (10 samples) below.

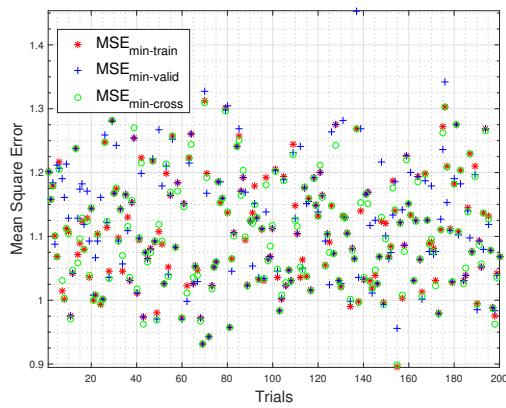


Figure 19: MSE on test data.

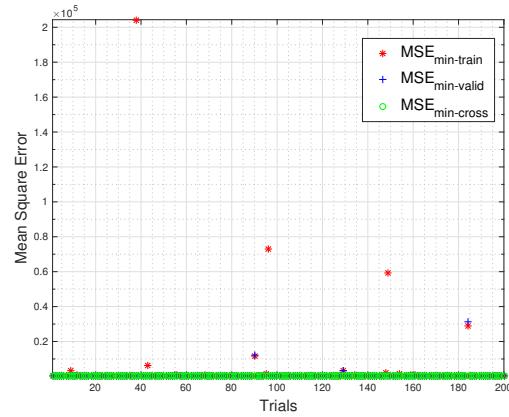


Figure 20: MSE on test data.

It can be seen by comparing the plots of Figure 19 and 20 that the range of values of the MSE vary by a substantial amount in Figure 20 which confirms the results obtained in Table 8 where all of the tuning methods with 10 training samples had large standard deviations indication large variations in the MSE values as shown in Figure 20.

Exercise 8

The MATLAB data file *boston.mat* was loaded into the workspace for processing as required as can be seen in code *ex9.m* which will be submitted with the assignment.

Exercise 9

(a)

Naive Regression was performed using vectors of ones for the training (\mathbf{x}_{train}) and test (\mathbf{x}_{test}) sets such that we would be learning the simple equation $\hat{y} = b$. The resultant MSE for the training and test sets can be seen in Table 10 shown below.

MSE Training	σ_{train}	MSE Test	σ_{test}
85.9461	4.8111	81.6110	9.4262

Table 10: Comparison of MSE performances with naive regression.

(b)

Based on the data loaded from the Boston housing data set. Linear regression was then performed on each one of the attributes while incorporating a bias term such that the weight vector $\mathbf{w} \in \mathbb{R}^2$ is learnt. The results obtained can be seen in Figures 21 - 23 below.

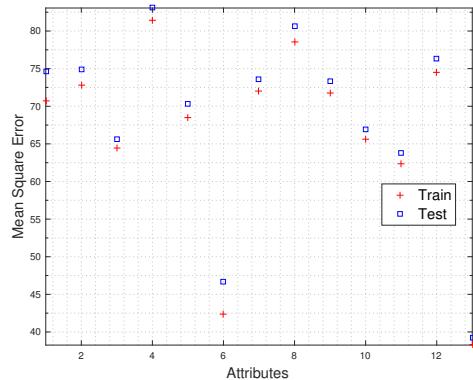


Figure 21: Average MSE for 13 attributes.

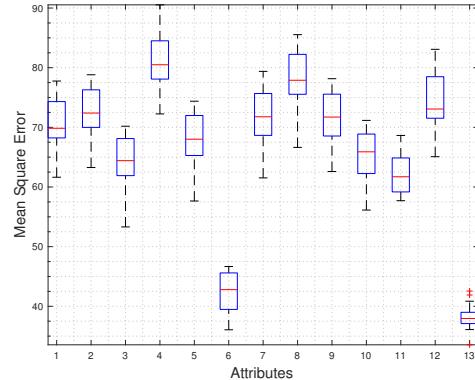


Figure 22: Box plot for MSE on the training set.

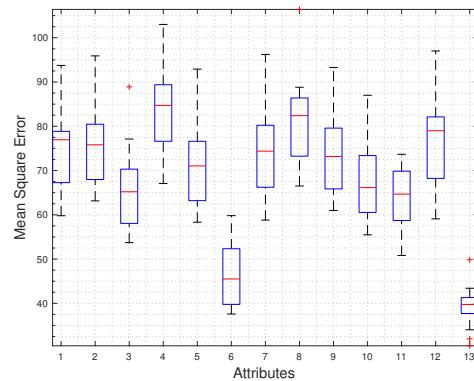


Figure 23: Box plot for MSE on the test set.

It can be seen from Figure 21 that the test errors were slightly larger than the training error. It can be seen that attributes 6 and 13 result in the lowest MSE for the training and test set (Figure 22).

23). Therefore, it can be concluded that they aid the most in predicting the output value $\hat{y} \approx y$.

The resultant mean and standard deviation for each attribute can be seen in Table 11 shown below.

Attribute No.	MSE Training	σ_{train}	MSE Test	σ_{test}
1	70.7447	4.2548	74.6865	9.1147
2	72.8396	3.9868	74.9556	8.2195
3	64.3904	4.4036	65.5785	8.9016
4	81.4900	4.7935	83.0593	9.4631
5	68.4933	4.2524	70.3485	8.5533
6	42.3054	3.5653	46.7082	7.2898
7	72.0200	4.5887	73.5604	9.3796
8	78.5866	4.6642	80.5866	9.5948
9	71.7204	4.2897	73.3460	8.6945
10	65.5571	4.0711	66.9019	8.2134
11	62.3123	3.6134	63.7735	7.3950
12	74.4986	4.5331	76.2912	9.2620
13	38.2555	2.0485	39.1809	4.2550

Table 11: Comparison of attribute MSE performances.

(c)

Linear regression was then performed using all the attributes available through the Boston housing data. The resultant training and test results can be seen in Table 12 below.

MSE Training	σ_{train}	MSE Test	σ_{test}
21.6457	1.2055	24.0536	2.9680

Table 12: Comparison of attribute MSE performances with all attributes.

It can be seen from Table 12 above that the resultant weight vector w performs far better than any of the individual regressors presented in Table 11 of Exercise 9(b). The figures demonstrating these results in a more graphical form can be seen in Figures 24 and 25 shown below.

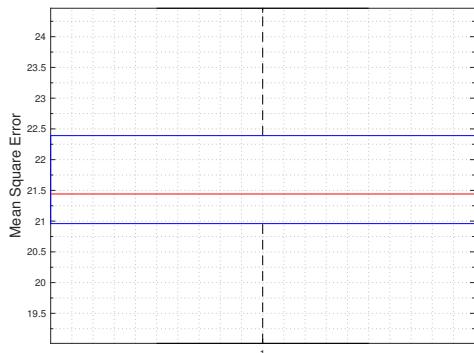


Figure 24: Box plot for MSE on the training set with all attributes.

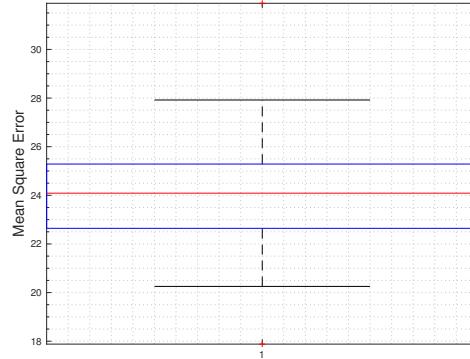


Figure 25: Box plot for MSE on the test set with all attributes.

Exercise 10

(a)

In this exercise Kernel Ridge Regression was computed on the Boston housing data in order to evaluate its performance in Exercise 10 against the performance of Linear Regression utilized in Exercise 9.

The function *kridgereg.m* was created with input parameters \mathbf{K} , \mathbf{y} and γ and output vector $\boldsymbol{\alpha}$ of dual weights. This function implements the following equation shown below.

$$\boldsymbol{\alpha} = (\mathbf{K} + \gamma l \mathbf{I}_l)^{-1} \mathbf{y} \quad (8)$$

The matlab code of the function *kridgereg.m* can be seen below.

```
1 function a = kridgereg(K, y, gamma)
2
3 a = (K + (gamma * size(y, 1) * eye(size(y, 1))) ) \ y;
```

In order to evaluate that the implementation is correct using the result from the assignment question page 6 the following equation was formed.

$$\boldsymbol{\alpha} = (\mathbf{X}\mathbf{X}' + \gamma l \mathbf{I}_l)^{-1} \mathbf{y} \quad (9)$$

This implementation correctly computes $\boldsymbol{\alpha}$ because the MATLAB left divide operation is used, instead of the explicit matrix inversion of the term $(\mathbf{K} + \gamma l \mathbf{I}_l)^{-1}$. The function correctly implements Equation 8 as required. Explicit matrix inversion is computationally expensive when compared to the left divide operation in MATLAB, thus it is computed slower. The residual error due to left divide operation is several orders of magnitude smaller than the explicit inverse function $\text{inv}((\mathbf{K} + \gamma l \mathbf{I}_l)^{-1})$ that could be used in MATLAB [2].

(b)

The function *dualcost.m* was created as required with the following input parameters \mathbf{K} , \mathbf{y} and the dual weight vector $\boldsymbol{\alpha}$ and returns the Mean Square Error (MSE) as shown below.

```
1 function mse = dualcost(K, y, a)
2
3 mse = (K*a - y)'*(K*a - y) / size(y, 1);
```

The function shown above implements Equation 10 shown below.

$$MSE = \frac{1}{l} (\mathbf{K}_{test} \boldsymbol{\alpha} - \mathbf{y})' (\mathbf{K}_{test} \boldsymbol{\alpha} - \mathbf{y}) \quad (10)$$

(c)

The vector of γ and σ were created as specified. Using the Gaussian Kernel as shown in Equation 11 shown below, Kernel Ridge Regression was performed.

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right) \quad (11)$$

(1)

Figure 26 shown below demonstrates the cross-validation error as a function of γ and σ .

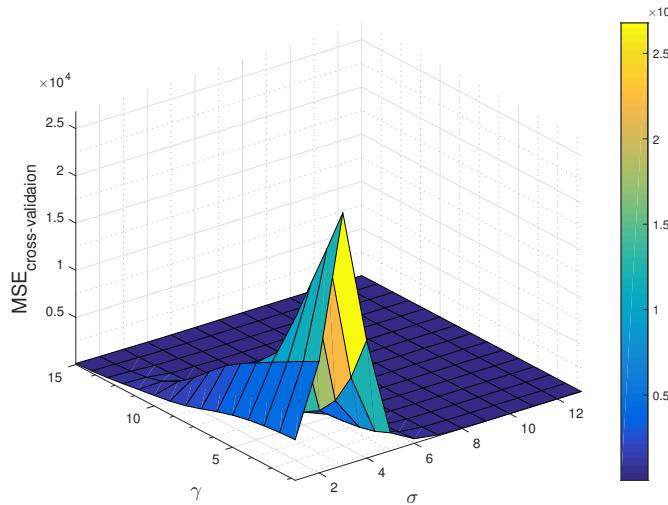


Figure 26: Cross validation error with Kernel Ridge Regression.

(2)

The optimal hyper-parameters γ^* and σ^* that minimized the training and test sets were found to be 1.8190×10^{-12} and 4096 respectively. The resultant MSE for the training and test sets can be seen in Table 13 below.

MSE train	MSE test
7.6349	14.4546

Table 13: MSE for optimal parameters γ^* and σ^* .

(d)

Exercise 10(a)-(c) was repeated over 20 random splits of the data and the resultant training and test errors were stored. Their relevant means and standard deviations were calculated and tabulated with the results obtained in Exercise 9 as can be seen in Table 14 shown below.

Method	MSE train	MSE test
Naive Regression	85.94 ± 4.81	81.61 ± 9.42
Linear Regression (attribute 1)	70.74 ± 4.25	74.68 ± 9.11
Linear Regression (attribute 2)	72.83 ± 3.98	74.95 ± 8.21
Linear Regression (attribute 3)	64.39 ± 4.40	65.57 ± 8.90
Linear Regression (attribute 4)	81.49 ± 4.79	83.05 ± 9.46
Linear Regression (attribute 5)	68.49 ± 4.25	70.34 ± 8.55
Linear Regression (attribute 6)	42.30 ± 3.56	46.70 ± 7.28
Linear Regression (attribute 7)	72.02 ± 4.58	73.56 ± 9.37
Linear Regression (attribute 8)	78.58 ± 4.66	80.58 ± 9.59
Linear Regression (attribute 9)	71.72 ± 4.28	73.34 ± 8.69
Linear Regression (attribute 10)	65.55 ± 4.07	66.90 ± 8.21
Linear Regression (attribute 11)	62.31 ± 3.61	63.77 ± 7.39
Linear Regression (attribute 12)	74.49 ± 4.53	76.29 ± 9.26
Linear Regression (attribute 13)	38.25 ± 2.04	39.18 ± 4.25
Kernel Ridge Regression	7.20 ± 1.53	19.89 ± 25.50

Table 14: Comparison of regression methods.

It can be seen from Table 14 that Kernel Ridge Regression performed the best when compared to other linear regressors. Thus, it is clearly evident from these results that performing linear regression on the non-linear dataset does not yield good predictive performance.

Therefore by using the Gaussian Kernel function corresponding to a non-linear feature map, linear regression using Ridge Regression can be utilised in the feature space, thus allowing the creation of a nonlinear regression function on the nonlinear dataset in the input space as stated in the paragraph above. This is not computational expensive as the Kernel trick makes use of inner products of the feature vector $K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)' \phi(\mathbf{x}_j)$ and rather than the actual vectors \mathbf{x}_i .

Question 1

In order to understand the relationship between linear regression with a Gaussian Kernel $K_\beta(\mathbf{x}, \mathbf{t}) = \exp(-\beta \|\mathbf{x} - \mathbf{t}\|^2)$ and 1-Nearest Neighbour classifier is to first understand the latter.

The 1-Nearest Neighbour algorithm can be broken down into three steps [3]:

1. Given training data $D = \{\mathbf{t}_i, y_i\}$ and an input \mathbf{x} data point.
2. Find index $j = \operatorname{argmin}_i d(\mathbf{x}, \mathbf{t}_i)$ that minimizes the euclidean distance ($\|\mathbf{x} - \mathbf{t}_i\|_2$).
3. $\hat{y} = \operatorname{sign}(y_j)$ where $\hat{y} \in \pm 1$. Therefore, the input \mathbf{x} has been classified as \hat{y} .

The downside of simple 1-Nearest Neighbour classification is that the nearest neighbour around a data point \mathbf{x} of interest will receive equal weight as can be seen in Step 3 of the algorithm shown above. The 1-NN algorithm can be adjusted to incorporate distance by means of the euclidean distance as can be seen in Equation 12 below.

$$\hat{y}(\mathbf{x}) = \operatorname{sign}(d(\mathbf{x}, \mathbf{t}_j)y_j) \quad (12)$$

$$\hat{y}(\mathbf{x}) = \operatorname{sign}(\|\mathbf{x} - \mathbf{t}_j\|_2 y_j) \quad (13)$$

However, when using the Gaussian Kernel in classification, the index number of the nearest neighbours ($K=1$) for an input \mathbf{x} does not need to be found, in order to classify. Instead, all neighbours surrounding \mathbf{x} can be summed up. Thus, data points (neighbours) closer to \mathbf{x} will receive a larger weight and contribute more to the summation as shown below.

$$\hat{y}(\mathbf{x}) = \operatorname{sign}\left(\sum_{j=1}^n K_\beta(\mathbf{x}, \mathbf{t}_j)y_j\right) \quad (14)$$

$$\hat{y}(\mathbf{x}) = \operatorname{sign}\left(\sum_{j=1}^n \exp\left(-\beta \|\mathbf{x} - \mathbf{t}_j\|^2\right)y_j\right) \quad (15)$$

$$\hat{y}(\mathbf{x}) = \operatorname{sign}\left(\sum_{j=1}^n \exp\left(-\frac{\|\mathbf{x} - \mathbf{t}_j\|^2}{2\sigma^2}\right)y_j\right) \quad (16)$$

It can be seen from Equation 15 and 16 take a similar form to that of Equation 13 shown earlier for 1-Nearest Neighbour. It is important to note now, that the parameter β or correspondingly so, σ defines the kernel width. In saying this, the parameter β affects how suddenly the weighting of surrounding data points \mathbf{t}_i reduces as their distance ($\|\mathbf{x} - \mathbf{t}_j\|^2$) increases.

The Gaussian Kernel in linear classification performs a similar weighting role as $d(\cdot, \cdot)$ mentioned earlier for 1-NN algorithm.

The variation of β is an important role to understand in Equation 15, thus we will consider two cases namely:

1. $\beta \rightarrow 0$: As β becomes smaller, the Gaussian Kernel will approach the value of 1 ($K_\beta(\mathbf{x}, \mathbf{t}_j) \rightarrow 1$), thus the weighting of each neighbour will become equally weighted regardless of distance. In effect, the width of the Gaussian Kernel is increasing, allowing more neighbours to weigh in.
2. $\beta \rightarrow \infty$: As β becomes larger, the width of the Gaussian Kernel becomes smaller ($K_\beta(\mathbf{x}, \mathbf{t}_j) \rightarrow 0$). Thus in effect only taking into account neighbours in the immediate vicinity surrounding \mathbf{x} . As β approaches zero, the classifier tends to the same behaviour as that of a 1-NN algorithm.

Question 2

(a)

In order to understand how we would change the perceptron to learn linear classifier with bias $b \Re$ as shown below in Equation 18.

$$f_{\mathbf{w},b}(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b) \quad (17)$$

We need to evaluate the perceptron to learn linear classifier with no bias $b = 0$ as shown in Equation 18 below.

$$f_{\mathbf{w},b}(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x}) \quad (18)$$

Using the notation defined in [1], let X denote the input space and Y denote the output domain. Where, $X \subseteq \Re^n$, $Y = \{-1, 1\}$. More formally, this can be labelled as the set training set S :

$$S = ((\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)) \subseteq (X \times Y)^m$$

The Perceptron algorithm pseudocode with no bias can be seen below in Algorithm 1 [1],[4]:

Algorithm 1 Perceptron Algorithm without bias

```

1: Given training set  $S$  and learning rate  $\eta \epsilon \Re^+$ 
2:  $\mathbf{w}_1 \leftarrow \mathbf{0}$ ,  $k \leftarrow 0$ 
3: for  $i = 1$  to  $m$ 
4:   if  $y_i(\text{sign}(\langle \mathbf{w}_k \cdot \mathbf{x}_i \rangle)) \leq 0$ 
5:      $\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k + \eta y_i \mathbf{x}_i$ 
6:      $k \leftarrow k + 1$ 
7:   end if
8: end for
9: return  $\mathbf{w}_k$ 
```

The Perceptron algorithm shown above was then updated to include the bias term b as can be seen in Algorithm 2 shown below.

Algorithm 2 Perceptron Algorithm with bias

```

1: Given training set  $S$  and learning rate  $\eta \epsilon \Re^+$ 
2:  $\mathbf{w}_1 \leftarrow \mathbf{0}$ ;  $k \leftarrow 0$ ;  $b_0 \leftarrow 0$ 
3:  $R \leftarrow \max_{1 \leq i \leq m} \|\mathbf{x}_i\|$ 
4: for  $i = 1$  to  $m$ 
5:   if  $y_i(\text{sign}(\langle \mathbf{w}_k \cdot \mathbf{x}_i \rangle + b_k)) \leq 0$ 
6:      $\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k + \eta y_i \mathbf{x}_i$ 
7:      $b_{k+1} \leftarrow b_k + \eta y_i R^2$ 
8:      $k \leftarrow k + 1$ 
9:   end if
10: end for
11: return  $\mathbf{w}_k, b_k$ 
```

It can be seen when comparing Algorithm 1 and 2 there only subtle differences the update procedure with b_k now included when a mistake is made, but the predominate procedure of the Perceptron remains unchanged.

Intuitively, the expression contained within the $\text{sign}(\cdot)$ function namely, $\langle \mathbf{w} \cdot \mathbf{x} \rangle + b$ splits the data S by means of a hyperplane that satisfies the equation:

$$\langle \mathbf{w} \cdot \mathbf{x} \rangle + b = 0$$

Thus, the separating of the data into two half spaces corresponding to the positive class $+1$ and the negative class -1 . The weight vector \mathbf{w} defines the direction perpendicular to the separating hyperplane. The term b is the bias term, which moves the hyperplane parallel to itself as can be seen in Figure 27 below.

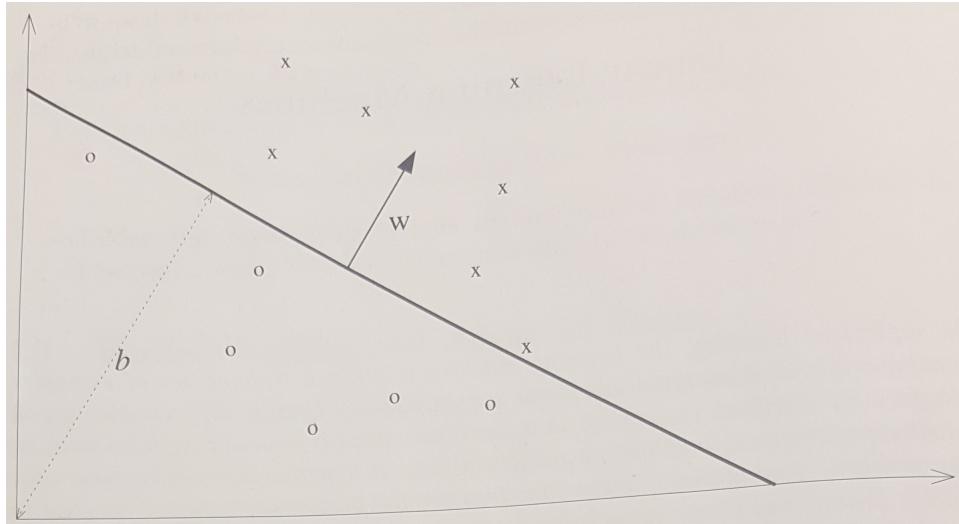


Figure 27: A separating hyperplane (\mathbf{w}, b) for a two dimensional set [1].

It is clear to see from the algorithms derived for the Perceptron and the geometric interpretation that by including a bias term b changes the threshold of the algorithm seen in line 5 of Algorithm 2 and update procedure. Thus, by varying b , it moves the hyperplane parallel to itself.

(b)

Firstly, we will consider the bound on the number of mistakes incurred by the Perceptron algorithm without bias. This result was derived through Novikoff Thereom [1].

$$k \leq \left(\frac{R_2}{\gamma} \right)^2 \quad (19)$$

Where:

- R_2 : $\|\mathbf{x}_i\| \leq R_2$ for all i .

We will now derive the upper bounds on the number of mistakes the perceptron makes when the bias term is incorporated.

Given a sequence of data:

$$S = [(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)]$$

Therefore, there must exist a weight vector \mathbf{w}^* such that is of unit length, $\|\mathbf{w}^*\| = 1$ and the following relationship hold true for all points in S .

$$y_i (\langle \mathbf{w}^* \cdot \mathbf{x}_i \rangle) \geq \gamma \quad (20)$$

We will derive the lower bound and the upper bounds on $\|\mathbf{w}^*\|$ to determine bound on perceptron mistakes. We will make a few changes to the \mathbf{x}_i and \mathbf{w} vectors to incorporate the bias term.

$$\begin{aligned}\hat{\mathbf{x}} &= (\mathbf{x}^T, R_1)^T \\ \hat{\mathbf{w}} &= (\mathbf{w}^T, \frac{b}{R_1})^T\end{aligned}$$

Where:

- $R_1: \|\hat{\mathbf{x}}_i\| \leq R_1$ for all i .

Therefore, let \mathbf{w}_{k-1} be the weight used, when the k-th mistake was made. Starting the algorithm with $\mathbf{w}_0 = \mathbf{0}$. Assuming on the k-th mistake using data point (\mathbf{x}_i, y_i) such that:

$$\text{sign}(\hat{\mathbf{w}}_{k-1}^T \hat{\mathbf{x}}_i) = \text{sign}(\mathbf{w}_{k-1}^T \mathbf{x}_i + b_{k-1}) \neq y_i$$

The incorrect classification implies:

$$y_i (\langle \hat{\mathbf{w}}_{k-1} \cdot \hat{\mathbf{x}}_i \rangle) = y_i (\langle \mathbf{w}_{k-1} \cdot \mathbf{x}_i \rangle + b_{k-1}) \leq 0 \quad (21)$$

Therefore, the Perceptron Algorithm would update in the following manner:

$$\hat{\mathbf{w}}_k = \hat{\mathbf{w}}_{k-1} + \eta y_i \hat{\mathbf{x}}_i = (\mathbf{w}_{k-1}^T, \frac{b_{k-1}}{R_1})^T + \eta y_i (\mathbf{x}_i^T, R_1)^T$$

Because the following is true [1]:

$$\frac{b_k}{R_1} = \frac{b_{k-1}}{R_1} + \eta y_i R_1$$

$$\langle \hat{\mathbf{w}}_k \cdot \hat{\mathbf{w}}^* \rangle = \langle \hat{\mathbf{w}}_{k-1} \cdot \hat{\mathbf{w}}^* \rangle + \eta y_i \langle \hat{\mathbf{x}}_i \cdot \hat{\mathbf{w}}^* \rangle \quad (22)$$

This equation above can be the lower bound using Equation 21. Thus, the following expression becomes evident:

$$\eta \gamma + \langle \hat{\mathbf{w}}_{k-1} \cdot \hat{\mathbf{w}}^* \rangle \leq \eta \gamma \langle \hat{\mathbf{x}}_i \cdot \hat{\mathbf{w}}^* \rangle + \langle \hat{\mathbf{w}}_{k-1} \cdot \hat{\mathbf{w}}^* \rangle$$

$$\eta \gamma + \langle \hat{\mathbf{w}}_{k-1} \cdot \hat{\mathbf{w}}^* \rangle \leq \langle \hat{\mathbf{w}}_k \cdot \hat{\mathbf{w}}^* \rangle$$

Therefore, by induction the following relationship is formed [5]:

$$k \eta \gamma \leq \langle \hat{\mathbf{w}}_k \cdot \hat{\mathbf{w}}^* \rangle \quad (23)$$

In a similar manner the upper bound can be formed:

$$\begin{aligned}\|\hat{\mathbf{w}}_k\|^2 &= \|\hat{\mathbf{w}}_{k-1} + \eta y_i \hat{\mathbf{x}}_i\|^2 \\ &= \|\hat{\mathbf{w}}_{k-1}\|^2 + 2\eta y_i \langle \hat{\mathbf{w}}_{k-1} \cdot \hat{\mathbf{x}}_i \rangle + \eta^2 \|\hat{\mathbf{x}}_i\|^2\end{aligned}$$

Using the fact that $y_i^2 = 1$ and Equation 20.

$$\|\hat{\mathbf{w}}_k\|^2 \leq \|\hat{\mathbf{w}}_{k-1}\|^2 + \eta^2 \|\hat{\mathbf{x}}_i\|^2$$

$$\|\hat{\mathbf{w}}_k\|^2 \leq \|\hat{\mathbf{w}}_{k-1}\|^2 + \eta^2 (\|\mathbf{x}_i\|^2 + R_1^2)$$

Using the fact that $\|\mathbf{x}_i\| \leq R_1$.

$$\|\mathbf{x}_i\|^2 \leq R_1^2$$

$$\|\hat{\mathbf{w}}_k\|^2 \leq \|\hat{\mathbf{w}}_{k-1}\|^2 + \eta^2 (R_1^2 + R_1^2)$$

$$\|\hat{\mathbf{w}}_k\|^2 \leq \|\hat{\mathbf{w}}_{k-1}\|^2 + 2\eta^2 R_1^2$$

Therefore, through induction the following inequality can be shown to exist:

$$\begin{aligned} \|\hat{\mathbf{w}}_k\|^2 &\leq 2\eta^2 k R_1^2 \\ \|\hat{\mathbf{w}}_k\| &\leq \sqrt{2k\eta} R_1 \end{aligned}$$

$$k\eta\gamma \leq \langle \hat{\mathbf{w}}_k \cdot \hat{\mathbf{w}}^* \rangle \leq \|\hat{\mathbf{w}}_k\| \|\hat{\mathbf{w}}^*\| \leq \sqrt{2k\eta} R_1 \|\hat{\mathbf{w}}^*\|$$

$$k\eta\gamma \leq \|\hat{\mathbf{w}}_k\| \|\hat{\mathbf{w}}^*\| \leq \sqrt{2k\eta} R_1 \|\hat{\mathbf{w}}^*\|$$

$$k\eta\gamma \leq \sqrt{2k\eta} R_1 \|\hat{\mathbf{w}}^*\|$$

$$\begin{aligned} k\eta &\leq \frac{\sqrt{2k\eta} R_1 \|\hat{\mathbf{w}}^*\|}{\gamma} \\ k &\leq \frac{2R_1^2 \|\hat{\mathbf{w}}^*\|^2}{\gamma} \end{aligned}$$

Therefore, due to the fact that $b^* \leq R_1$ must hold true for a non-trivial separation of the data [1], and the fact that $\|\mathbf{w}^*\| = 1$ [6].

$$\|\hat{\mathbf{w}}^*\|^2 \leq \|\mathbf{w}^*\|^2 + 1 = 1 + 1 = 2$$

$$\begin{aligned} k &\leq \frac{4R_1^2}{\gamma^2} \\ k &\leq \left(\frac{2R_1}{\gamma} \right)^2 \end{aligned} \tag{24}$$

It is important to note the relationship R_1 and R_2 have in determining whether the bounds increases, assuming γ remains unchanged. It can be shown at $R_1 > R_2$ in the following manner:

$$R_1 > R_2$$

$$\|\hat{\mathbf{x}}_i\| > \|\mathbf{x}_i\|$$

$$\sqrt{1^2 + \sum_{i=1}^m x_i^2} > \sqrt{\sum_{i=1}^m x_i^2}$$

Therefore when comparing Equations 19 and 24, it is clear that the bound does increase by a factor of 4. Thus, the bound on the number of mistakes increases by adding a bias term in the Perceptron Algorithm.

Question 3

(a)

K_c is positive semidefinite kernel for $c \geq 0$.

$$K_c(x, z) := c + \sum_{i=1}^n x_i z_i \quad (25)$$

Since $\sum_{i=1}^n x_i z_i$ is a linear kernel it is positive semidefinite.

By the property we know that sum of two positive semidefinites is a positive semidefinite. Hence if we assume c is 1×1 matrix and if it satisfies $y C y' \geq 0 \forall y \in \mathbb{R}$ condition, then we can say c is positive semidefinite.

We also know that every 1×1 dimension matrix is positive semidefinite if $c_{11} \geq 0$. Therefore, for K_c to be a positive semidefinite, we need $c \geq 0$

(b)

The c parameter does not directly influence kernel but controls the trade-off between the margin and the size of the slack variables. [1]

Question 4

(a)

The four algorithms were considered in order to evaluate the sample complexity of each, namely *The Perceptron*, *Winnow*, *Least Squares* and *1-nearest neighbour* [4].

In this question, we are concerned with looking at the relationship between the (n) number of features against the number of training examples (m) required which would lead to a reduction in the generalisation error. More formally, the "working definition" as stated in the question page, "...sample complexity is the minimum number of examples (m) to incur no more than 10% generalisation error (on average)" [7].

Thus, the algorithms mentioned earlier were, implemented and their performance was evaluated. The sample complexity plots for each algorithm can be seen in Figures 28 - 31 shown below.

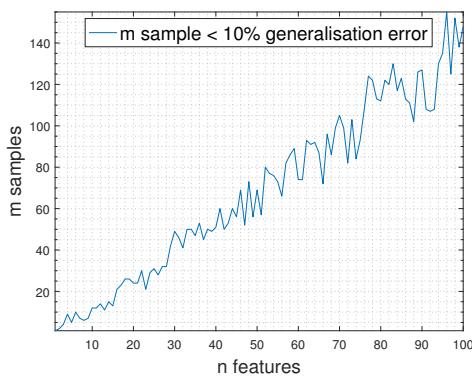


Figure 28: Estimated sample complexity of *Perceptron* algorithm.

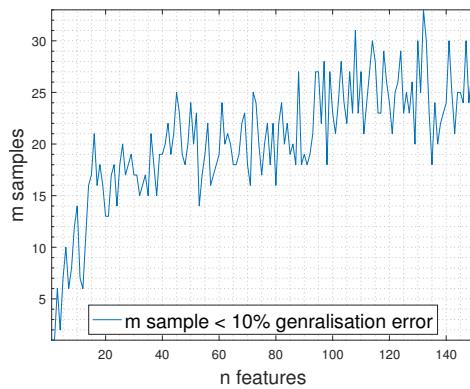


Figure 29: Estimated sample complexity of *Winnow* algorithm.

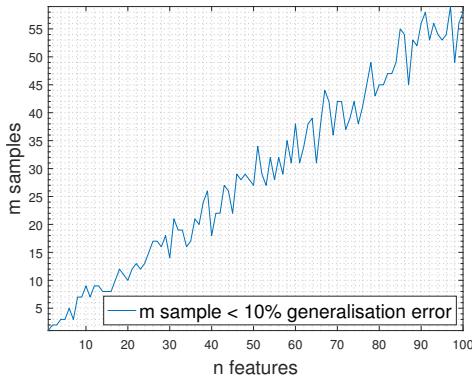


Figure 30: Estimated sample complexity of *Least Squares* algorithm.

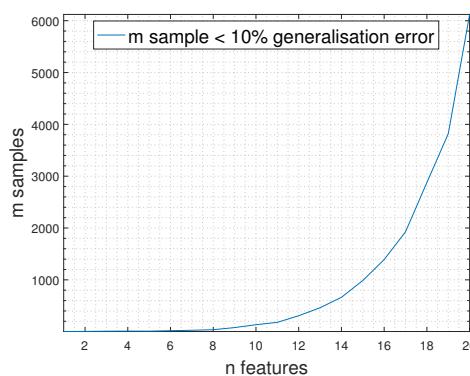


Figure 31: Estimated sample complexity of *1-Nearest Neighbour* algorithm.

It can be seen from Figure 28 - 31 shown above that for each algorithm as the number of features increase, the corresponding number of training samples increase. Therefore, the sample complexity increases as a function of features.

All the MATLAB code required to generate these results can be seen in Appendix .

(b)

It became evident through empirically testing in Q4(a), that computing the "sample complexity" exactly through simulation became extremely expensive.

In order to test the sample complexity exactly, we iteratively increased n dimensions of the features of the input data. And for each such vector, we iteratively trained and tested on those trained parameters by varying the lengths of samples m . For each combination of n feature and m samples, we computed the percent error of mis-classification. For each value of n we evaluated the percent error and continued to test, until the error for a particular value of m is not at least 10%. When error drops below or equal to 10%, we record the value of m and do the same process for the next value of n .

The computational time for each algorithm to converge to a solution as the number of samples (m) and number of features (n) were increased can be seen in the following table.

Algorithm	Mean Comp. Time	$\sigma_{comptime}$	n	m	Test sample size
Perceptron	337s	$\pm 10s$	100	500	500
Winnow	198s	$\pm 6s$	100	500	500
Least Squares	140s	$\pm 3s$	100	500	500
1-NN	21682s	$\pm 200s$	20	10000	500

Table 15: Computational time for algorithms.

It can be seen from Table 15 that *1 Nearest Neighbour* algorithm took the longest, whereas *Least Square* algorithm was the fastest amongst all.

Thus, in order to evaluate the sample complexity for each algorithm certain trade-off were made namely:

1. Number of independent trial were reduced: In order to find the generalized error, we need to take test on multiple trials and then average out the errors. We had to reduce them to 100 and trade off on the accuracy for speed of execution.
2. Dimension (n) was limited for each algorithm: In case of *1-NN* we had to limit the dimension to just 20 features which took around 6 hours to compute. Whereas we limited the dimension n for other algorithms at 100.
3. In order to faster execute the sample complexity over the range, we further enhanced the method to iterate over odd values of the dimensions (n) and samples (m), at the cost of accuracy.

(c)

In order to fully understand the relationship of m vs. n , the Figures 28 - 31 generated in Q4(a) were analysed closer.

For Perceptron algorithm [8] it can be seen from Figure 28 that as $n \rightarrow \infty$, the sample complexity increased as a linear function of dimension ($m = \Theta(n)$). It is found that Perceptron make mistakes in each trial and is $\Omega(kN)$ mistakes.

If we then consider the complexity of the Winnow algorithm, it can be deduced from Figure 29 that the sample complexity is a logarithmic function of dimension ($m = \Theta(\log(n))$). As we increase the dimensions (n) the number of samples (m) increase at a certain extent and then slows down and nearly stays constant at ∞ . The mistake bound is $\Omega(kln2^n)$

In both Perceptron algorithm and Winnow, the new hypothesis depends on the old hypothesis at a learning rate $\eta > 0$.

The sampling complexity of *Least Squares* algorithm was found to be linear in dimensions (n) with respect to sampling (m). This algorithm is the fastest to run out of other three. The complexity was found to be $m = \Theta(n)$ for $n \rightarrow \infty$ from Figure 30.

Finally, the 1-Nearest Neighbour algorithm was inspected and the found that the sample complexity is exponential in nature with respect to dimensions ($n \rightarrow \infty$) from 31. 1-NN performance was seen to be worst out of all four algorithms. It took more than 6 hours to run for just up to 20 dimensions. It made most errors and hence as a result it requires more number of samples to converge for generalization error below 10%

Appendices

Exercise 1: Code

In this and all of the following sections of the code we have used external functions at multiple places. Please refer Library Functions section at the end of the appendix for implementation.

ex1.m

```
1 %-----%
2 % Module: GI01 – Supervised Learning
3 % Assignment : Coursework 1
4 % Author : Russel Daries , Nitish Mutha
5 % Student ID: 16079408, 15113106
6 % Question: 1
7 % Section: Part 1
8 % Description: Least Squares Regression
9 %-----%
10
11 % clearing memory
12 clear all
13 close all
14 clc
15
16 addpath('.. / library')
17
18 %% Ex1(a)
19
20 w = randn(1,1);
21 x = randn(600,1);
22 n = randn(600,1);
23
24 %% Ex1(b)
25 part_b_result = calculateLSR(x,w,n,100,500)
26
27
28 %% Ex1(c)
29 part_c_result = calculateLSR(x,w,n,10,500)
30
31 %% Ex1(d)
32
33 result_100 = [];
34 result_10 = [];
35 for j = 1:200
36     w_j = randn(1,1);
37     x_j = randn(600,1);
38     n_j = randn(600,1);
39     result_100 = [result_100 ; calculateLSR(x_j,w_j,n_j,100,500)];
40     result_10 = [result_10 ; calculateLSR(x_j,w_j,n_j,10,500)];
41 end
42
43 average_result_100 = mean(result_100)
```

```
44 average_result_10 = mean(result_10)
45
46 % Code to plot resultant figures
47
48 figure;
49 plot(result_10(:,2), 'k*')
50 hold on
51 plot(result_100(:,2), 'b+')
52 hold on
53 plot(average_result_10(2)*ones(length(result_10),1), 'r—')
54 hold on
55 plot(average_result_100(2)*ones(length(result_100),1), 'm—')
56 xlabel('Trial number', 'FontSize', 15)
57 ylabel('Mean Squared Error', 'FontSize', 15)
58 grid on;
59 set(gcf, 'Color', 'w');
60 leg=legend('train10_{err}', 'train100_{err}', 'train-err10_{avg}', 'train-err100_{avg}', 'Location', 'Best')
61 set(leg, 'FontSize', 15)
62 set(gca, 'YMinorTick', 'on')
63 grid minor
64 axis tight;
65 print('ex1d_err_train', '-depsc')
66 close all;
67
68
69 figure
70 plot(result_10(:,3), 'k*')
71 grid on;
72 hold on
73 plot(result_100(:,3), 'b+')
74 hold on
75 plot(average_result_10(2)*ones(length(result_10),1), 'c—')
76 hold on
77 plot(average_result_100(3)*ones(length(result_100),1), 'r—')
78 hold on
79 xlabel('Trial number', 'FontSize', 15)
80 ylabel('Mean Squared Error', 'FontSize', 15)
81 set(gcf, 'Color', 'w');
82 leg=legend('test10_{err}', 'test100_{avg}', 'test10-err_{avg}', 'test100-err_{avg}', 'Location', 'Best');
83 set(leg, 'FontSize', 15)
84 set(gca, 'YMinorTick', 'on')
85 grid minor
86 axis tight;
87 ylim([0.8 2.2])
88 print('ex1d_err_test', '-depsc')
89 close all;
90
```

```
91 figure;
92 plot(result_10(:,1), 'k*')
93 hold on
94 plot(result_100(:,1), 'b+')
95 hold on
96 plot(average_result_10(1)*ones(length(result_10),1), 'r—')
97 hold on
98 plot(average_result_100(1)*ones(length(result_100),1), 'm—')
99 xlabel('Trial number', 'FontSize', 15)
100 ylabel('Weight Estimate (w)', 'FontSize', 15)
101 grid on;
102 set(gcf, 'Color', 'w');
103 leg=legend('w-train10', 'w-train100', 'w-train-err10_{avg}', 'w-train-err100_{avg}
    }', 'Location', 'Best')
104 set(leg, 'FontSize', 15)
105 set(gca, 'YMinorTick', 'on')
106 grid minor
107 axis tight;
108 print('ex1d_err_wtrain', '-depsc')
109 close all;
```

Exercise 2: Code**ex2.m**

```
1 %-----%
2 % Course : Supervised Learning
3 % Assignment : 1
4 % Excercise 2
5 % Programmers: Russel Daries , Nitish Mutha
6 % Student ID: 16079408, 15113106
7 % Question: 2
8 % Section: Part 1
9 % Description: Least Squares Regression
10 %-----%
11
12 clear all
13 close all
14 clc
15
16 addpath('.. / library')
17
18 %% Ex2(a)
19 w = randn(10,1);
20 x = randn(600,10);
21 n = randn(600,1);
22
23 %% Ex2(b)
24 [ part_b_result_w , part_b_result_mse_train , part_b_result_mse_test ] =
25 calculateLSRex2(x,w,n,100,500);
26 [ part_c_result_w , part_c_result_mse_train , part_c_result_mse_test ] =
27 calculateLSRex2(x,w,n,10,500);
28
29 % Initilizae variables
30 result_100 = [];
31 result_10 = [];
32 w_100 = [];
33 w_10 = [];
34 mse_train_100 = [];
35 mse_train_10 = [];
36 mse_test_100 = [];
37 mse_test_10 = [];
38
39 for j = 1:200
40     % Generate dataset
41     w_j = randn(10,1);
42     x_j = randn(600,10);
43     n_j = randn(600,1);
44
45     % Store resultant calculation
46     % 100 training examples
```

```
45 [ result_w_100_j , result_mse_train_100_j , result_mse_test_100_j ] =
46     calculateLSRex2(x_j , w_j , n_j ,100 ,500) ;
47 w_100 = [ w_100; result_w_100_j ];
48 mse_train_100 = [ mse_train_100 ; result_mse_train_100_j ];
49 mse_test_100 = [ mse_test_100; result_mse_test_100_j ];
50
51 % 10 training example
52 [ result_w_10_j , result_mse_train_10_j , result_mse_test_10_j ] =
53     calculateLSRex2(x_j , w_j , n_j ,10 ,500) ;
54 w_10 = [ w_10; result_w_10_j ];
55 mse_train_10 = [ mse_train_10 ; result_mse_train_10_j ];
56 mse_test_10 = [ mse_test_10; result_mse_test_10_j ];
57
58 % Average MSE for 10,100 training examples
59 average_mse_train_100 = mean( mse_train_100 )
60 average_mse_test_100 = mean( mse_test_100 )
61
62 % Average MSE on the test sets
63 average_mse_train_10 = mean( mse_train_10 )
64 average_mse_test_10 = mean( mse_test_10 )
65
66 % Code to plot resultant figures
67
68 figure;
69 h1 = plot( mse_train_10 , 'k*' )
70 hold on
71 h2 = plot( mse_train_100 , 'b+' )
72 hold on
73 h3 = plot( average_mse_train_10*ones( length( mse_train_10 ) ,1) , 'r—' )
74 hold on
75 h4 = plot( average_mse_train_100*ones( length( mse_test_100 ) ,1) , 'm—' )
76 xlabel( 'Trial number' , 'FontSize' ,15)
77 ylabel( 'Mean Squared Error' , 'FontSize' ,15)
78 grid on;
79 set(gcf , 'Color' , 'w');
80 leg=legend( 'train10_{err}' , 'train100_{err}' , 'train-err10_{avg}' , 'train-err100_{avg}' , 'Location' , 'Best' )
81 set(leg , 'FontSize' ,15)
82 set(gca , 'YMinorTick' , 'on')
83 grid minor
84 axis tight;
85 print( 'ex2b_err_train' , '-depsc' )
86 close all;
87
88 figure
89 plot( mse_test_10 , 'k*' )
```

```
91 grid on;
92 hold on
93 plot(mse_test_100, 'b+')
94 hold on
95 plot(average_mse_test_10*ones(length(mse_test_10),1), 'm—')
96 hold on
97 plot(average_mse_test_100*ones(length(mse_test_100),1), 'r—')
98 hold on
99 xlabel('Trial number', 'FontSize', 15)
100 ylabel('Mean Squared Error', 'FontSize', 15)
101 set(gcf, 'Color', 'w');
102 leg=legend('test10_{err}', 'test100_{avg}', 'test10-err_{avg}', 'test100-err_{avg}
    }', 'Location', 'Best');
103 set(leg, 'FontSize', 15)
104 set(gca, 'YMinorTick', 'on')
105 grid minor
106 axis tight;
107 ylim([0.0 3e3])
108 print('ex2b_err_test', '-depsc')
109 close all;
110
111 figure;
112 plot(w_10, 'k*')
113 hold on
114 plot(w_100, 'b+')
115 hold on
116 plot(mean(w_10)*ones(length(w_10),1), 'r—')
117 hold on
118 plot(mean(w_100)*ones(length(w_100),1), 'm—')
119 xlabel('Trial number', 'FontSize', 15)
120 ylabel('Weight Estimate (w)', 'FontSize', 15)
121 grid on;
122 set(gcf, 'Color', 'w');
123 leg=legend('w-train10', 'w-train100', 'w-train10_{avg}', 'w-train100_{avg}', '
    Location', 'Best');
124 set(leg, 'FontSize', 15)
125 set(gca, 'YMinorTick', 'on')
126 grid minor
127 axis tight;
128 print('ex2b_wtrain', '-depsc')
129 close all;
```

Exercise 4: Code**ex4.m**

```
1 %-----%
2 % Module: GI01 – Supervised Learning
3 % Assignment : Coursework 1
4 % Author : Russel Stuart Daries , Nitish Mutha
5 % Student ID: 16079408 ,15113106
6 % Question: 4
7 % Section: Part 1
8 % Description: Effect of Regularisation Parameter
9 %-----%
10
11 %%%
12 clear all
13 close all
14 clc
15
16 addpath('.. / library')
17
18 %% Ex4(a)
19 w = randn(10,1);
20 x = randn(600,10);
21 n = randn(600,1);
22 gamma_power = -6:1:3;
23 % Range of gamma values
24 gamma = 10.^gamma_power;
25
26 %Initialise variables
27 result_mse_train_100 = [];
28 result_mse_test_100 = [];
29 for i = 1 : size(gamma,2)
30     [w_estimator_100 , mse_train_100 , mse_test_100] = calculateLSReg4(x,w,n
31                 ,100,500,gamma(i));
32     result_mse_train_100 = [result_mse_train_100 ; mse_train_100];
33     result_mse_test_100 = [result_mse_test_100 ; mse_test_100];
34 end
35 figure
36 plot(log(gamma) ,result_mse_train_100 , 'b-o');
37 hold on;
38 plot(log(gamma) ,result_mse_test_100 , 'r-*');
39 grid on;
40 set(gcf , 'Color' , 'w');
41 xlabel('{log(\gamma)}' , 'FontSize' ,15);
42 ylabel('Mean Square Error' , 'FontSize' ,15);
43 leg=legend('MSE-train_{100}' , 'MSE-test_{100}' , 'Location' , 'Best');
44 set(leg , 'FontSize' ,15);
45 set(gca , 'YMinorTick' , 'on');
```

```
46 grid minor
47 axis tight;
48 print('ex4a_mse_100','-depsc');
49 close all;
50
51
52 %% Ex4(b)
53 result_mse_train_10 = [];
54 result_mse_test_10 = [];
55 for i = 1 : size(gamma,2)
56 [w_estimator_10 , mse_train_10 , mse_test_10] = calculateLSRex4(x,w,n
57 ,10,500,gamma(i));
58 result_mse_train_10 = [result_mse_train_10 ; mse_train_10];
59 result_mse_test_10 = [result_mse_test_10 ; mse_test_10];
60 end
61 figure
62 plot(log(gamma) ,result_mse_train_10 , 'b-o');
63 hold on;
64 plot(log(gamma) ,result_mse_test_10 , 'r-*');
65 grid on;
66 set(gcf, 'Color', 'w');
67 xlabel('{log(\gamma)}', 'FontSize',15);
68 ylabel('Mean Square Error', 'FontSize',15);
69 leg=legend('MSE-train_{10}', 'MSE-test_{10}', 'Location', 'Best');
70 set(leg, 'FontSize',15);
71 set(gca, 'YMinorTick', 'on');
72 grid minor
73 axis tight;
74 print('ex4b_mse_10','-depsc');
75 close all;
76
77
78 %% Ex4(c)
79
80 result_mse_train_100 = [];
81 result_mse_train_10 = [];
82 result_mse_test_100 = [];
83 result_mse_test_10 = [];
84 result_w_100 = [];
85 result_w_10 = [];
86
87 % 200 independant trials
88 for j = 1:200
89
90 w_j = randn(10,1);
91 x_j = randn(600,10);
92 n_j = randn(600,1);
93 w_100 = [];
```

```

94     w_10 = [];
95
96     mse_train_100 = [];
97     mse_train_10 = [];
98     mse_test_100 = [];
99     mse_test_10 = [];
100
101    % Varying regularisation parameter
102    for i=1: size(gamma,2)
103        [result_w_100_j , result_mse_train_100_j , result_mse_test_100_j] =
104            calculateLSRex4(x_j,w_j,n_j,100,500,gamma(i));
105        w_100 = [w_100 , result_w_100_j];
106        mse_train_100 = [mse_train_100 , result_mse_train_100_j];
107        mse_test_100 = [mse_test_100 , result_mse_test_100_j];
108
109        [result_w_10_j , result_mse_train_10_j , result_mse_test_10_j] =
110            calculateLSRex4(x_j,w_j,n_j,10,500,gamma(i));
111        w_10 = [w_10 , result_w_10_j];
112        mse_train_10 = [mse_train_10 , result_mse_train_10_j];
113        mse_test_10 = [mse_test_10 , result_mse_test_10_j];
114
115    end
116
117    result_mse_train_100 = [result_mse_train_100 ; mse_train_100];
118    result_mse_train_10 = [result_mse_train_10 ; mse_train_10];
119
120    result_mse_test_100 = [result_mse_test_100 ; mse_test_100];
121    result_mse_test_10 = [result_mse_test_10 ; mse_test_10];
122
123 end
124
125 % Computing the average result
126 average_mse_train_100 = mean(result_mse_train_100);
127 average_mse_test_100 = mean(result_mse_test_100);
128 average_mse_train_10 = mean(result_mse_train_10);
129 average_mse_test_10 = mean(result_mse_test_10);
130
131
132 figure
133 plot(log(gamma) ,average_mse_train_100 , 'b-*');
134 hold on;
135 grid on;
136 plot(log(gamma) ,average_mse_test_100 , 'r-o');
137 set(gcf , 'Color' , 'w');
138 xlabel('log(\gamma)' , 'FontSize' ,15);
139 ylabel('Mean Square Error' , 'FontSize' ,15);
140 leg=legend('MSE-train_{100}' , 'MSE-test_{100}' , 'Location' , 'Best');

```

```
141 set(leg , 'FontSize' ,15);
142 set(gca , 'YMinorTick' , 'on' );
143 grid minor
144 axis tight;
145 print('ex4c_mse_200_100' , '-depsc');
146 close all;
147
148 figure;
149 plot(log(gamma) ,average_mse_train_10 , 'b-*');
150 hold on; grid on;
151 plot(log(gamma) ,average_mse_test_10 , 'r-o');
152 grid on;
153 xlabel('{log(\gamma)}' , 'FontSize' ,15);
154 ylabel('Mean Square Error' , 'FontSize' ,15);
155 leg=legend('MSE-train_{10}' , 'MSE-test_{10}' , 'Location' , 'Best');
156 set(leg , 'FontSize' ,15);
157 set(gca , 'YMinorTick' , 'on' );
158 grid minor
159 axis tight;
160 print('ex4c_mse_200_10' , '-depsc');
161 close all;
```

Exercise 5: Code**ex5.m**

```
1 %-----%
2 % Module: GI01 – Supervised Learning
3 % Assignment : Coursework 1
4 % Author : Russel Stuart Daries , Nitish Mutha
5 % Student ID: 16079408 ,15113106
6 % Question: 5
7 % Section: Part 1
8 % Description: Tuning Regularization Parameter with Validation Set
9 %-----%
10
11 %%%
12 clear all
13 close all
14 clc
15
16 addpath('.. / library')
17
18 %% Ex5(a) and Ex5(b)
19
20 gamma_power = -6:1:3;
21 gamma = 10.^gamma_power;
22
23 result_mse_train_100 = [];
24 result_mse_train_10 = [];
25
26 result_mse_valid_100 = [];
27 result_mse_valid_10 = [];
28
29 result_mse_test_100 = [];
30 result_mse_test_10 = [];
31
32 result_w_100 = [];
33 result_w_10 = [];
34
35 best_gamma_100 = zeros(200,1);
36 best_gamma_10 = zeros(200,1);
37
38 g_average_mse_test_100 = [];
39 g_average_mse_test_10 = [];
40
41 for j = 1:200
42
43     w_j = randn(10,1);
44     x_j = randn(600,10);
45     n_j = randn(600,1);
46     w_100 = [];
```

```

47 w_10 = [];
48
49 mse_train_100 = [];
50 mse_valid_100 = [];
51 mse_test_100 = [];
52
53 mse_train_10 = [];
54 mse_valid_10 = [];
55 mse_test_10 = [];
56
57
58 for i=1: size(gamma,2)
    % Calculate MSE for each data set with 100 training samples
    [result_w_100_j , result_mse_train_100_j , result_mse_valid_100_j ,
     result_mse_test_100_j] = calculateLSRex5(x_j , w_j , n_j ,80 ,20 ,500 ,
     gamma(i));
    w_100 = [w_100; result_w_100_j];
    mse_train_100 = [mse_train_100 , result_mse_train_100_j];
    mse_valid_100 = [mse_valid_100 , result_mse_valid_100_j];
    mse_test_100 = [mse_test_100 , result_mse_test_100_j];
    % Calculate MSE for each data set with 10 training samples
    [result_w_10_j , result_mse_train_10_j , result_mse_valid_10_j ,
     result_mse_test_10_j] = calculateLSRex5(x_j , w_j , n_j ,8 ,2 ,500 ,gamma(
     i));
    w_10 = [w_10; result_w_10_j];
    mse_train_10 = [mse_train_10 , result_mse_train_10_j];
    mse_valid_10 = [mse_valid_10 , result_mse_valid_10_j];
    mse_test_10 = [mse_test_10 , result_mse_test_10_j];
71
72 end
73 % Find optimal gamma value from validation set MSE
74 [min_mse_100 , idx_min_100] = min(mse_valid_100);
75 best_gamma_100(j ,1) = gamma(idx_min_100);
76
77 [min_mse_10 , idx_min_10] = min(mse_valid_10);
78 best_gamma_10(j ,1) = gamma(idx_min_10);
79
80 result_w_100 = [result_w_100 ; w_100];
81 result_w_10 = [result_w_10 ; w_10];
82
83 result_mse_train_100 = [result_mse_train_100 ; mse_train_100];
84 result_mse_train_10 = [result_mse_train_10 ; mse_train_10];
85
86 result_mse_valid_100 = [result_mse_valid_100 ; mse_valid_100];
87 result_mse_valid_10 = [result_mse_valid_10 ; mse_valid_10];
88
89 result_mse_test_100 = [result_mse_test_100 ; mse_test_100];
90 result_mse_test_10 = [result_mse_test_10 ; mse_test_10];
91

```

```
92 [ w_test_100 , g_result_mse_test_100 ] = calculateLSRTTest( x_j , w_j , n_j  
93 , 80 , 20 , 500 , gamma( idx_min_100 ) );  
94 g_average_mse_test_100 = [ g_average_mse_test_100 , g_result_mse_test_100 ];  
95  
96 [ w_test_10 , g_result_mse_test_10 ] = calculateLSRTTest( x_j , w_j , n_j , 8 , 2 , 500 ,  
97 gamma( idx_min_10 ) );  
98 g_average_mse_test_10 = [ g_average_mse_test_10 , g_result_mse_test_10 ];  
99  
100  
101 % Compute averages  
102 average_mse_train_100 = mean( result_mse_train_100 );  
103 average_mse_valid_100 = mean( result_mse_valid_100 );  
104 average_mse_test_100 = mean( result_mse_test_100 );  
105  
106 average_mse_train_10 = mean( result_mse_train_10 );  
107 average_mse_valid_10 = mean( result_mse_valid_10 );  
108 average_mse_test_10 = mean( result_mse_test_10 );  
109  
110 % plots for mse vs gamma values for 100  
111 figure;  
112 plot( log( gamma ) , average_mse_train_100 , 'r-*' );  
113 hold on;  
114 plot( log( gamma ) , average_mse_valid_100 , 'g-o' );  
115 plot( log( gamma ) , average_mse_test_100 , 'b+-' );  
116 grid on;  
117 set( gcf , 'Color' , 'w' );  
118 xlabel( 'Regularization parameter (\gamma)' , 'FontSize' , 15 );  
119 ylabel( 'Mean Square Error' , 'FontSize' , 15 );  
120 leg=legend( 'MSE_{train}' , 'MSE_{validation}' , 'MSE_{test}' , 'Location' , 'Best' );  
121 set( leg , 'FontSize' , 15 );  
122 set( gca , 'YMinorTick' , 'on' );  
123 grid minor  
124 axis tight;  
125 print( 'ex5a_80_20' , '-depsc' );  
126 close all;  
127  
128 % plots for mse vs gamma values for 10  
129 figure;  
130 plot( log( gamma ) , average_mse_train_10 , 'r-*' );  
131 hold on;  
132 plot( log( gamma ) , average_mse_valid_10 , 'g-o' );  
133 plot( log( gamma ) , average_mse_test_10 , 'b+-' );  
134 grid on;  
135 set( gcf , 'Color' , 'w' );  
136 xlabel( 'Regularization parameter (\gamma)' , 'FontSize' , 15 );  
137 ylabel( 'Mean Square Error' , 'FontSize' , 15 );  
138 leg=legend( 'MSE_{train}' , 'MSE_{validation}' , 'MSE_{test}' , 'Location' , 'Best' );
```

```
139 set(leg , 'FontSize' ,15);
140 set(gca , 'YMinorTick' , 'on');
141 grid minor
142 axis tight;
143 print('ex5b_8_2' , '-depsc');
144 close all;
145
146 % plots for test mse vs experiments for 100 and 10
147 figure;
148 plot(g_average_mse_test_100 , 'r*');
149 hold on;
150 plot(g_average_mse_test_10 , 'b+');
151 grid on;
152 set(gcf , 'Color' , 'w');
153 xlabel('Experiments' , 'FontSize' ,15);
154 ylabel('Mean Square Error' , 'FontSize' ,15);
155 leg=legend( 'MSE_{100}' , 'MSE_{10}' , 'Location' , 'Best');
156 set(leg , 'FontSize' ,15);
157 set(gca , 'YMinorTick' , 'on');
158 grid minor
159 axis tight;
160 ylim([0.0 40])
161 print('ex5b_mse_optimal_gamma' , '-depsc');
162 close all;
163
164 %% Ex5(c)
165 mean_gamma_100 = mean(best_gamma_100)
166 mean_gamma_10 = mean(best_gamma_10)
167
168 %% Ex5(d)
169
170 % Repeating tuning of regularization parameter with validation set with
171 % data similair to Ex1(a)-(d)
172 result_mse_train_100 = [];
173 result_mse_train_10 = [];
174
175 result_mse_valid_100 = [];
176 result_mse_valid_10 = [];
177
178 result_mse_test_100 = [];
179 result_mse_test_10 = [];
180
181 result_w_100 = [];
182 result_w_10 = [];
183
184 best_gamma_100 = zeros(200,1);
185 best_gamma_10 = zeros(200,1);
186
187 g_average_mse_test_100 = [];
```

```

188 g_average_mse_test_10 = [];
189
190 for j = 1:200
191
192 w_j = randn(1,1);
193 x_j = randn(600,1);
194 n_j = randn(600,1);
195 w_100 = [];
196 w_10 = [];
197
198 mse_train_100 = [];
199 mse_valid_100 = [];
200 mse_test_100 = [];
201
202 mse_train_10 = [];
203 mse_valid_10 = [];
204 mse_test_10 = [];
205
206
207 for i=1: size(gamma,2)
208 % Calculate MSE for each data set with 100 training samples
209 [result_w_100_j, result_mse_train_100_j, result_mse_valid_100_j,
210     result_mse_test_100_j] = calculateLSRex5(x_j, w_j, n_j, 80, 20, 500,
211     gamma(i));
212 w_100 = [w_100; result_w_100_j];
213 mse_train_100 = [mse_train_100, result_mse_train_100_j];
214 mse_valid_100 = [mse_valid_100, result_mse_valid_100_j];
215 mse_test_100 = [mse_test_100, result_mse_test_100_j];
216 % Calculate MSE for each data set with 10 training samples
217 [result_w_10_j, result_mse_train_10_j, result_mse_valid_10_j,
218     result_mse_test_10_j] = calculateLSRex5(x_j, w_j, n_j, 8, 2, 500, gamma(
219     i));
220 w_10 = [w_10; result_w_10_j];
221 mse_train_10 = [mse_train_10, result_mse_train_10_j];
222 mse_valid_10 = [mse_valid_10, result_mse_valid_10_j];
223 mse_test_10 = [mse_test_10, result_mse_test_10_j];
224
225
226 end
227 % Find optimal gamma value from validation set MSE
228 [min_mse_100, idx_min_100] = min(mse_valid_100);
229 best_gamma_100(j,1) = gamma(idx_min_100);
230
231 [min_mse_10, idx_min_10] = min(mse_valid_10);
232 best_gamma_10(j,1) = gamma(idx_min_10);
233
234
235 result_w_100 = [result_w_100 ; w_100];
236 result_w_10 = [result_w_10 ; w_10];
237
238 result_mse_train_100 = [result_mse_train_100 ; mse_train_100];

```

```
233     result_mse_train_10 = [result_mse_train_10 ; mse_train_10];  
234  
235     result_mse_valid_100 = [result_mse_valid_100 ; mse_valid_100];  
236     result_mse_valid_10 = [result_mse_valid_10 ; mse_valid_10];  
237  
238     result_mse_test_100 = [result_mse_test_100 ; mse_test_100];  
239     result_mse_test_10 = [result_mse_test_10 ; mse_test_10];  
240  
241     [w_test_100 , g_result_mse_test_100] = calculateLSRTTest(x_j,w_j,n_j  
242     ,80,20,500,gamma(idx_min_100));  
243     g_average_mse_test_100 = [g_average_mse_test_100 , g_result_mse_test_100];  
244  
245     [w_test_10 , g_result_mse_test_10] = calculateLSRTTest(x_j,w_j,n_j,8,2,500,  
246     gamma(idx_min_10));  
247     g_average_mse_test_10 = [g_average_mse_test_10 , g_result_mse_test_10];  
248  
249  
250 % Compute averages  
251 average_mse_train_100 = mean(result_mse_train_100);  
252 average_mse_valid_100 = mean(result_mse_valid_100);  
253 average_mse_test_100 = mean(result_mse_test_100);  
254  
255 average_mse_train_10 = mean(result_mse_train_10);  
256 average_mse_valid_10 = mean(result_mse_valid_10);  
257 average_mse_test_10 = mean(result_mse_test_10);  
258  
259 % plots for mse vs gamma values for 100  
260 figure;  
261 plot(log(gamma),average_mse_train_100,'r-*');  
262 hold on;  
263 plot(log(gamma),average_mse_valid_100,'g-o');  
264 plot(log(gamma),average_mse_test_100,'b+-');  
265 grid on;  
266 set(gcf, 'Color', 'w');  
267 xlabel('Regularization parameter (\gamma)', 'FontSize', 15);  
268 ylabel('Mean Square Error', 'FontSize', 15);  
269 leg=legend('MSE-{train}', 'MSE-{validation}', 'MSE-{test}', 'Location', 'Best');  
270 set(leg, 'FontSize', 15);  
271 set(gca, 'YMinorTick', 'on');  
272 grid minor  
273 axis tight;  
274 print('ex5d_80_20',' -depsc');  
275 close all;  
276  
277 % plots for mse vs gamma values for 10  
278 figure;  
279 plot(log(gamma),average_mse_train_10,'r-*');
```

```
280 hold on;
281 plot(log(gamma),average_mse_valid_10,'g-o');
282 plot(log(gamma),average_mse_test_10,'b+-');
283 grid on;
284 set(gcf, 'Color', 'w');
285 xlabel('Regularization parameter (\gamma)', 'FontSize', 15);
286 ylabel('Mean Square Error', 'FontSize', 15);
287 leg=legend('MSE_{train}', 'MSE_{validation}', 'MSE_{test}', 'Location', 'Best');
288 set(leg, 'FontSize', 15);
289 set(gca, 'YMinorTick', 'on');
290 grid minor;
291 axis tight;
292 print('ex5d_8_2', '-depsc');
293 close all;
```

Exercise 6: Code**ex6.m**

```
1 %-----%
2 % Module: GI01 – Supervised Learning
3 % Assignment : Coursework 1
4 % Author : Russel Stuart Daries , Nitish Mutha
5 % Student ID: 16079408 ,15113106
6 % Question: 6
7 % Section: Part 1
8 % Description: Tuning Regularization Parameter with Cross–Validation
9 %-----%
10
11 %%%
12 clear all
13 close all
14 clc
15
16 addpath('.. / library')
17
18 %% Ex6(a) and Ex6(b)
19
20 gamma_power = -6:1:3;
21 gamma = 10.^gamma_power;
22
23 result_mse_train_100 = [];
24 result_mse_train_10 = [];
25
26 result_mse_valid_100 = [];
27 result_mse_valid_10 = [];
28
29 result_mse_test_100 = [];
30 result_mse_test_10 = [];
31
32 result_w_100 = [];
33 result_w_10 = [];
34
35
36 g_average_mse_test_100 = [];
37 g_average_mse_test_10 = [];
38
39 w_j = randn(10,1);
40 x_j = randn(600,10);
41 n_j = randn(600,1);
42
43 % Perform 5 fold cross validation
44 for j = 1:5
45
46
```

```
47 w_100 = [];
48 w_10 = [];

49
50 mse_train_100 = [];
51 mse_valid_100 = [];

52
53
54 mse_train_10 = [];
55 mse_valid_10 = [];

56
57
58 for i=1: size(gamma,2)
59     [ result_w_100_j , result_mse_train_100_j , result_mse_valid_100_j ] =
60         calculateLSRex6(x_j,w_j,n_j,100, gamma(i),j);
61     w_100 = [w_100; result_w_100_j];
62     mse_train_100 = [mse_train_100 , result_mse_train_100_j];
63     mse_valid_100 = [mse_valid_100 , result_mse_valid_100_j];

64     [ result_w_10_j , result_mse_train_10_j , result_mse_valid_10_j ] =
65         calculateLSRex6(x_j,w_j,n_j,10, gamma(i),j);
66     w_10 = [w_10; result_w_10_j];
67     mse_train_10 = [mse_train_10 , result_mse_train_10_j];
68     mse_valid_10 = [mse_valid_10 , result_mse_valid_10_j];

69
70 end
71 %Identify gamma that minimises MSE
72 [min_mse_100 , idx_min_100] = min(mse_valid_100);
73 best_gamma_100(j,1) = gamma(idx_min_100);

74 [min_mse_10 , idx_min_10] = min(mse_valid_10);
75 best_gamma_10(j,1) = gamma(idx_min_10);

76
77 result_w_100 = [result_w_100 ; w_100];
78 result_w_10 = [result_w_10 ; w_10];

79
80 result_mse_train_100 = [result_mse_train_100 ; mse_train_100];
81 result_mse_train_10 = [result_mse_train_10 ; mse_train_10];

82
83 result_mse_valid_100 = [result_mse_valid_100 ; mse_valid_100];
84 result_mse_valid_10 = [result_mse_valid_10 ; mse_valid_10];

85
86 end
87
88 mean_gamma_100 = mean(best_gamma_100)
89 mean_gamma_10 = mean(best_gamma_10)

90
91 % Average all dataset errors
92 average_mse_train_100 = mean(result_mse_train_100);
93 average_mse_valid_100 = mean(result_mse_valid_100);
```

```
94
95 average_mse_train_10 = mean(result_mse_train_10);
96 average_mse_valid_10 = mean(result_mse_valid_10);
97
98 mse_test_100 = [];
99 mse_test_10 = [];
100 % Calculate MSE for test set
101 for i=1: size(gamma,2)
102     [result_w_100_j , result_mse_train_100_j , result_mse_test_100_j] =
103         calculateLSRex4(x_j , w_j , n_j ,100 ,500 ,gamma(i));
104     mse_test_100 = [mse_test_100 , result_mse_test_100_j];
105
106     [result_w_10_j , result_mse_train_10_j , result_mse_test_10_j] =
107         calculateLSRex4(x_j , w_j , n_j ,10 ,500 ,gamma(i));
108     mse_test_10 = [mse_test_10 , result_mse_test_10_j];
109
110 end
111
112 % plots for mse vs gamma values for 100
113 figure;
114 plot(log(gamma) ,average_mse_train_100 , 'r-*');
115 hold on;
116 plot(log(gamma) ,average_mse_valid_100 , 'b-o');
117 plot(log(gamma) ,mse_test_100 , 'g+-');
118 grid on;
119 set(gcf , 'Color' , 'w');
120 xlabel('log (\gamma)' , 'FontSize' ,15);
121 ylabel('Mean Square Error' , 'FontSize' ,15);
122 leg=legend( 'MSE_{avg train}' , 'MSE_{avg cross-val}' , 'MSE_{test}' , 'Location' ,
123             'Best');
124 set(leg , 'FontSize' ,15);
125 set(gca , 'YMinorTick' , 'on');
126 grid minor
127 axis tight;
128 print('ex6_cv_100' , '-depsc');
129 close all;
130
131 % plots for mse vs gamma values for 10
132 figure;
133 plot(log(gamma) ,average_mse_train_10 , 'r-*');
134 hold on;
135 plot(log(gamma) ,average_mse_valid_10 , 'b-o');
136 plot(log(gamma) ,mse_test_10 , 'g+-');
137 grid on;
138 set(gcf , 'Color' , 'w');
139 xlabel('log (\gamma)' , 'FontSize' ,12);
140 ylabel('Mean Square Error' , 'FontSize' ,12);
141 leg=legend( 'MSE_{avg train}' , 'MSE_{avg cross-val}' , 'MSE_{test}' , 'Location' ,
142             'Best');
```

```
139 set(leg , 'FontSize' ,15) ;  
140 set(gca , 'YMinorTick' , 'on') ;  
141 grid minor  
142 axis tight;  
143 print('ex6_cv_10' , '-depsc') ;  
144 close all;
```

Exercise 7: Code**ex7.m**

```
1 %-----%
2 % Module: GI01 – Supervised Learning
3 % Assignment : Coursework 1
4 % Author : Russel Stuart Daries , Nitish Mutha
5 % Student ID: 16079408 ,15113106
6 % Question: 7
7 % Section: Part 1
8 % Description: Comparing tuning–methods for regularisation parameters
9 %-----%
10
11 %%%
12 clear all
13 close all
14 clc
15
16 addpath('.. / library')
17
18 gamma_power = -6:1:3;
19 gamma = 10.^gamma_power;
20
21 %% Initialize variables of interest
22 test_4_100 = zeros(200,1);
23 test_4_10 = zeros(200,1);
24
25 test_5_100 = zeros(200,1);
26 test_5_10 = zeros(200,1);
27
28 test_6_100 = zeros(200,1);
29 test_6_10 = zeros(200,1);
30
31 all_gamma_4_100 = zeros(200,1);
32 all_gamma_4_10 = zeros(200,1);
33
34 all_gamma_5_100 = zeros(200,1);
35 all_gamma_5_10 = zeros(200,1);
36
37 all_gamma_6_100 = zeros(200,1);
38 all_gamma_6_10 = zeros(200,1);
39
40 %% Complete 200 independent trials
41 for j = 1:200
42
43     w_j = randn(10,1);
44     x_j = randn(600,10);
45     n_j = randn(600,1);
46
```

```
47 %% Select best gamma by minimizing training error
48 [gamma_4_100, gamma_4_10, mse_test_4_100, mse_test_4_10] = ex7_4(x_j, w_j, n_j,
   gamma);
49 %% Select best gamma by minimizing validation set error
50 [gamma_5_100, gamma_5_10, mse_test_5_100, mse_test_5_10] = ex7_5(x_j, w_j, n_j,
   gamma);
51 %% Select best gamma using 5-fold cross-validation
52 [gamma_6_100, gamma_6_10, mse_test_6_100, mse_test_6_10] = ex7_6(x_j, w_j, n_j,
   gamma);

53
54 % Assigning resultant mse test values to their respective outputs
55 test_4_100(j) = mse_test_4_100;
56 test_4_10(j) = mse_test_4_10;

57
58 test_5_100(j) = mse_test_5_100;
59 test_5_10(j) = mse_test_5_10;

60
61 test_6_100(j) = mse_test_6_100;
62 test_6_10(j) = mse_test_6_10;

63
64 % Assigning optimal gamma values to vectors
65 all_gamma_4_100(j) = gamma_4_100;
66 all_gamma_4_10(j) = gamma_4_10;

67
68 all_gamma_5_100(j) = gamma_5_100;
69 all_gamma_5_10(j) = gamma_5_10;

70
71 all_gamma_6_100(j) = gamma_6_100;
72 all_gamma_6_10(j) = gamma_6_10;

73 end

74
75 %% Average MSE values
76 avg_test_4_100 = mean(test_4_100)
77 avg_test_4_10 = mean(test_4_10)

78
79 avg_test_5_100 = mean(test_5_100)
80 avg_test_5_10 = mean(test_5_10)

81
82 avg_test_6_100 = mean(test_6_100)
83 avg_test_6_10 = mean(test_6_10)

84
85 %% Average gamma values
86 avg_gamma_4_100 = mean(all_gamma_4_100)
87 avg_gamma_4_10 = mean(all_gamma_4_10)

88
89 avg_gamma_5_100 = mean(all_gamma_5_100)
90 avg_gamma_5_10 = mean(all_gamma_5_10)

91
92 avg_gamma_6_100 = mean(all_gamma_6_100)
```

```
93 avg_gamma_6_10 = mean(all_gamma_6_10)
94
95 %% Standard deviation calculation
96 sd_4_100 = standard_deviation(avg_test_4_100,test_4_100)
97 sd_4_10 = standard_deviation(avg_test_4_10,test_4_10)
98
99 sd_5_100 = standard_deviation(avg_test_5_100,test_5_100)
100 sd_5_10 = standard_deviation(avg_test_5_10,test_5_10)
101
102 sd_6_100 = standard_deviation(avg_test_6_100,test_6_100)
103 sd_6_10 = standard_deviation(avg_test_6_10,test_6_10)
104
105 %% plots
106 figure;
107 plot(test_4_100,'r*');
108 hold on;
109 plot(test_5_100,'b+')
110 hold on
111 plot(test_6_100,'go');
112 grid on;
113 set(gcf,'Color','w');
114 xlabel('Trials','FontSize',15);
115 ylabel('Mean Square Error','FontSize',15);
116 leg=legend('MSE_{min-train}','MSE_{min-valid}','MSE_{min-cross}','Location','Best');
117 set(leg,'FontSize',15);
118 set(gca,'YMinorTick','on');
119 grid minor
120 axis tight;
121 print('ex7_mse_compare_tuning_100','-depsc');
122 close all;
123
124 figure;
125 plot(test_4_10,'r*');
126 hold on;
127 plot(test_5_10,'b+')
128 hold on
129 plot(test_6_10,'go');
130 grid on;
131 set(gcf,'Color','w');
132 xlabel('Trials','FontSize',15);
133 ylabel('Mean Square Error','FontSize',15);
134 leg=legend('MSE_{min-train}','MSE_{min-valid}','MSE_{min-cross}','Location','Best');
135 set(leg,'FontSize',15);
136 set(gca,'YMinorTick','on');
137 grid minor
138 axis tight;
139 print('ex7_mse_compare_tuning_10','-depsc');
```

```
140 close all;
141
142 figure;
143 plot(all_gamma_4_10, 'r*');
144 hold on;
145 plot(all_gamma_5_10, 'b+')
146 hold on
147 plot(all_gamma_6_10, 'go');
148 grid on;
149 set(gcf, 'Color', 'w');
150 xlabel('Trials', 'FontSize', 15);
151 ylabel('Mean Square Error', 'FontSize', 15);
152 leg=legend('MSE_{min-train}', 'MSE_{min-valid}', 'MSE_{min-cross}', 'Location', 'Best');
153 set(leg, 'FontSize', 15);
154 set(gca, 'YMinorTick', 'on');
155 grid minor
156 axis tight;
157 print('ex7_mse_compare_tuning_gamma_10', '-depsc');
158 close all;
159
160 figure;
161 plot(all_gamma_4_100, 'r*');
162 hold on;
163 plot(all_gamma_5_100, 'b+')
164 hold on
165 plot(all_gamma_6_100, 'go');
166 grid on;
167 set(gcf, 'Color', 'w');
168 xlabel('Trials', 'FontSize', 15);
169 ylabel('Mean Square Error', 'FontSize', 15);
170 leg=legend('MSE_{min-train}', 'MSE_{min-valid}', 'MSE_{min-cross}', 'Location', 'Best');
171 set(leg, 'FontSize', 15);
172 set(gca, 'YMinorTick', 'on');
173 grid minor
174 axis tight;
175 print('ex7_mse_compare_tuning_gamma_100', '-depsc');
176 close all;
```

ex7_4.m

```
1 function [ best_gamma_100, best_gamma_10, result_mse_test_100, result_mse_test_10 ] =
2     = ex7_4(x, w, n, gamma)
3
4 addpath('../library')
5 mse_train_100 = [];
6 mse_train_10 = [];
7 % Evaluate performance of test set
```

```

8 for i=1: size(gamma,2)
9 [ result_w_100_j , result_mse_train_100_j , result_mse_test_100_j ] =
   calculateLSRex4(x,w,n,100,500,gamma(i));
10 mse_train_100 = [ mse_train_100 , result_mse_train_100_j ];
11
12 [ result_w_10_j , result_mse_train_10_j , result_mse_test_10_j ] =
   calculateLSRex4(x,w,n,10,500,gamma(i));
13 mse_train_10 = [ mse_train_10 , result_mse_train_10_j ];
14
15 end
16 % Find the gamma that minimizes the MSE of training set
17 [ min_train_error_100 , min_idx_100 ] = min( mse_train_100 );
18 best_gamma_100 = gamma( min_idx_100 );
19
20 [ min_train_error_10 , min_idx_10 ] = min( mse_train_10 );
21 best_gamma_10 = gamma( min_idx_10 );
22
23 % Evaluate performance on test vector
24 [ result_w_100 , result_mse_train_100 , result_mse_test_100 ] = calculateLSRex4(x,
   w,n,100,500,best_gamma_100);
25 [ result_w_10 , result_mse_train_10 , result_mse_test_10 ] = calculateLSRex4(x,w,n
   ,10,500,best_gamma_10);

```

ex7_5.m

```

1 function [ best_gamma_100 , best_gamma_10 , result_mse_test_100 , result_mse_test_10 ]
   = ex7_5(x,w,n,gamma)
2
3 addpath( '../library' )
4
5 mse_valid_100 = [];
6 mse_valid_10 = [];
7 % Evaluate performance of training and validation sets of data
8 for i=1: size(gamma,2)
9   [ result_w_100_j , result_mse_train_100_j , result_mse_valid_100_j ,
   result_mse_test_100_j ] = calculateLSRex5(x,w,n,80,20,500,gamma(i));
10  mse_valid_100 = [ mse_valid_100 , result_mse_valid_100_j ];
11
12 [ result_w_10_j , result_mse_train_10_j , result_mse_valid_10_j ,
   result_mse_test_10_j ] = calculateLSRex5(x,w,n,8,2,500,gamma(i));
13  mse_valid_10 = [ mse_valid_10 , result_mse_valid_10_j ];
14
15 end
16 % Find the gamma that minimizes the MSE of validation set
17 [ min_mse_100 , idx_min_100 ] = min( mse_valid_100 );
18 best_gamma_100 = gamma( idx_min_100 );
19
20 [ min_mse_10 , idx_min_10 ] = min( mse_valid_10 );
21 best_gamma_10 = gamma( idx_min_10 );
22

```

```

23 % Evaluate performance on test vector
24 [ w_test_100 , result_mse_test_100 ] = calculateLSRTTest(x,w,n,80,20,500,
    best_gamma_100);
25 [ w_test_10 , result_mse_test_10 ] = calculateLSRTTest(x,w,n,8,2,500, best_gamma_10
    );

```

ex7_6.m

```

1 function [ best_gamma_100 , best_gamma_10 , result_mse_test_100 , result_mse_test_10 ]
    = ex7_6(x,w,n,gamma)
2
3 addpath( '../library' )
4
5 result_mse_valid_100 = [];
6 result_mse_valid_10 = [];
7
8 % Perform 5-fold cross-validation
9 for j = 1:5
10
11     mse_valid_100 = [];
12     mse_valid_10 = [];
13
14     for i=1: size(gamma,2)
15         [ result_w_100_j , result_mse_train_100_j , result_mse_valid_100_j ] =
16             calculateLSRex6(x,w,n,100, gamma(i),j);
17         mse_valid_100 = [ mse_valid_100 , result_mse_valid_100_j ];
18
19         [ result_w_10_j , result_mse_train_10_j , result_mse_valid_10_j ] =
20             calculateLSRex6(x,w,n,10, gamma(i),j);
21         mse_valid_10 = [ mse_valid_10 , result_mse_valid_10_j ];
22     end
23
24     result_mse_valid_100 = [ result_mse_valid_100 ; mse_valid_100 ];
25     result_mse_valid_10 = [ result_mse_valid_10 ; mse_valid_10 ];
26 end
27
28 % Find the gamma that minimizes the mean cross-validation MSE
29 [ min_mse_100 , min_idx_100 ] = min(mean(result_mse_valid_100));
30 best_gamma_100 = gamma(min_idx_100);
31
32 [ min_mse_10 , min_idx_10 ] = min(mean(result_mse_valid_10));
33 best_gamma_10 = gamma(min_idx_10);
34
35 % Evaluate performance on test vector
36 [ result_w_100 , result_mse_train_100 , result_mse_test_100 ] = calculateLSRex4(x,w
    ,n,100,500, best_gamma_100);
37 [ result_w_10 , result_mse_train_10 , result_mse_test_10 ] = calculateLSRex4(x,w,n
    ,10,500, best_gamma_10);

```

Exercise 9: Code**Ex9_a.m**

```
1 %-----%
2 % Module: GI01 – Supervised Learning
3 % Assignment : Coursework 1
4 % Author : Russel Stuart Daries , Nitish Mutha
5 % Student ID: 16079408 ,15113106
6 % Question: 9a
7 % Section: Part 1
8 % Description: Baseline versus full linear regression .
9 %-----%
10
11 close all;
12 clear all;
13 addpath('..//library')
14
15 load('boston.mat')
16
17 %% Ex9(a)
18
19 % 20 trials
20 for j=1:20
21
22     data_train = datasample(boston,round(2/3*size(boston,1)),1,'Replace',false);
23     % Stored left over rows in the test set 1/3
24     data_test = setdiff(boston, data_train, 'rows');
25
26     x_train_v = ones(round(2/3*size(boston,1)),1);
27     y_train = data_train(:,14);
28
29     x_test_v = ones(round(1/3*size(boston,1)),1);
30     y_test = data_test(:,14);
31
32     % Calculate weight vector w
33     w_star = learnModel(x_train_v, y_train);
34     % Calculate MSE
35     mse_train(j) = meanSquareError(w_star, x_train_v, y_train, size(x_train_v,1));
36     ;
37     mse_test(j) = meanSquareError(w_star, x_test_v, y_test, size(x_test_v,1));
38
39 end
40 % Calculate mean and standard deviation of trian and test set.
41 avg_mse_train = mean(mse_train)
42 std_train = std(mse_train)
43 avg_mse_test = mean(mse_test)
44 std_test = std(mse_test)
```

```

45
46 mse = [ avg_mse_train; avg_mse_test ]';

```

Ex9_b.m

```

1 %-----%
2 % Module: GI01 – Supervised Learning
3 % Assignment : Coursework 1
4 % Author : Russel Stuart Daries , Nitish Mutha
5 % Student ID: 16079408 ,15113106
6 % Question: 9b
7 % Section: Part 1
8 % Description: Baseline versus full linear regression .
9 %-----%
10
11 close all;
12 clear all;
13 addpath('..//library')
14
15 load('boston.mat')
16
17 %% Ex 9(b)
18
19 %% 20 trials
20 for j=1:20
21
22     % Radomly withdraw 2/3 data from Boston for training samples
23     data_train = datasample(boston,round(2/3*size(boston,1)),1,'Replace',false);
24     % Stored left over rows in the test set 1/3
25     data_test = setdiff(boston, data_train, 'rows');
26
27     % Segment the columns for training attributes
28     x_train = data_train(:,1:13);
29     % Take the last column as the true value we are trying to predict
30     y_train = data_train(:,14);
31     % Perform same mapping for test set
32     x_test = data_test(:,1:13);
33     y_test = data_test(:,14);
34
35     % Loop through each individual attribute with bias and try to predict
36     % outcome
37     for i=1:size(x_train,2)
38         x_train_v = [x_train(:,i), ones(size(x_train,1),1)];
39         w_star = learnModel(x_train_v, y_train);
40         mse_train(j,i) = meanSquareError(w_star, x_train_v, y_train, size(
41             x_train_v,1));
42         x_test_v = [x_test(:,i), ones(size(x_test,1),1)];

```

```
43     mse_test(j,i) = meanSquareError(w_star,x_test_v,y_test,size(x_test_v
44     ,1));
45
46 end
47
48 % Calculate mean and standard deviation for each of the 20 trials for the
49 % 13 attributes for train and test set
50 avg_mse_train = mean(mse_train)
51 std_train = std(mse_train)
52
53 avg_mse_test = mean(mse_test)
54 std_test = std(mse_test)
55
56 mse = [avg_mse_train; avg_mse_test]';
57
58 %% plots
59 figure;
60 plot(avg_mse_train,'r+');
61 hold on
62 plot(avg_mse_test,'bs');
63 set(gcf, 'Color', 'w');
64 xlabel('Attributes','FontSize',15);
65 ylabel('Mean Square Error','FontSize',15);
66 leg=legend('Train','Test','Location','Best');
67 set(leg,'FontSize',15);
68 set(gca,'YMinorTick','on');
69 grid minor
70 axis tight;
71 print('ex9b_avg_mse',' -depsc');
72 close all;
73
74 figure;
75 boxplot(mse_train);
76 set(gcf, 'Color', 'w');
77 xlabel('Attributes','FontSize',15);
78 ylabel('Mean Square Error','FontSize',15);
79 leg=legend('Train','Test','Location','Best');
80 set(leg,'FontSize',15);
81 set(gca,'YMinorTick','on');
82 grid minor
83 axis tight;
84 print('ex9b_mse_box_train',' -depsc');
85 close all;
86
87 figure;
88 boxplot(mse_test);
89 set(gcf, 'Color', 'w');
90 xlabel('Attributes','FontSize',15);
```

```

91 ylabel('Mean Square Error', 'FontSize',15);
92 leg=legend('Train','Test','Location','Best');
93 set(leg,'FontSize',15);
94 set(gca,'YMinorTick','on');
95 grid minor;
96 axis tight;
97 print('ex9b_mse_box_test','-depsc');
98 close all;

```

Ex9_c.m

```

1 % %
2 % Module: GI01 – Supervised Learning %
3 % Assignment : Coursework 1 %
4 % Author : Russel Stuart Daries , Nitish Mutha %
5 % Student ID: 16079408 ,15113106 %
6 % Question: 9c %
7 % Section: Part 1 %
8 % Description: Baseline versus full linear regression . %
9 %

10
11
12 close all;
13 clear all;
14 addpath('.. / library')
15
16 load('boston.mat')
17
18 %% Ex9(c)
19
20 % 20 trials
21 for j=1:20
22
23     % Radomly withdraw 2/3 data from Boston for training samples
24     data_train = datasample(boston,round(2/3*size(boston,1)),1,'Replace',false);
25
26     % Stored left over rows in the test set 1/3
27     data_test = setdiff(boston, data_train, 'rows');
28
29     % Segment the columns for training attributes
30     x_train = data_train(:,1:13);
31
32     % Take the last column as the true value we are trying to predict
33     y_train = data_train(:,14);
34
35     x_test = data_test(:,1:13);
36     y_test = data_test(:,14);
37
38     % Train the model vector w
39     x_train_v = [x_train, ones(size(x_train,1),1)];
40     w_star = learnModel(x_train_v, y_train);

```

```
38     mse_train(j) = meanSquareError(w_star, x_train_v, y_train, size(x_train_v,1))
39     ;
40
41     % Test the learnt vector w
42     x_test_v = [x_test, ones(size(x_test,1),1)];
43     mse_test(j) = meanSquareError(w_star, x_test_v, y_test, size(x_test_v,1));
44
45 end
46 avg_mse_train = mean(mse_train)
47 std_train = std(mse_train)
48
49 avg_mse_test = mean(mse_test)
50 std_test = std(mse_test)
51
52 mse = [avg_mse_train; avg_mse_test]';
53
54 %% plots
55 figure;
56 boxplot(mse_train);
57 set(gcf, 'Color', 'w');
58 % xlabel('Attributes', 'FontSize', 15);
59 ylabel('Mean Square Error', 'FontSize', 15);
60 leg=legend('Train', 'Test', 'Location', 'Best');
61 set(leg, 'FontSize', 15);
62 set(gca, 'YMinorTick', 'on');
63 grid minor
64 axis tight;
65 print('ex9c_mse_box_train', '-depsc');
66 close all;
67
68 figure;
69 boxplot(mse_test);
70 set(gcf, 'Color', 'w');
71 % xlabel('Attributes', 'FontSize', 15);
72 ylabel('Mean Square Error', 'FontSize', 15);
73 leg=legend('Train', 'Test', 'Location', 'Best');
74 set(leg, 'FontSize', 15);
75 set(gca, 'YMinorTick', 'on');
76 grid minor
77 axis tight;
78 print('ex9c_mse_box_test', '-depsc');
79 close all;
```

Exercise 10: Code

Ex10_RR.m

```
1 %-----%
2 % Module: GI01 – Supervised Learning
3 % Assignment : Coursework 1
4 % Author : Russel Stuart Daries , Nitish Mutha
5 % Student ID: 16079408 ,15113106
6 % Question: 10
7 % Section: Part 1
8 % Description: Kernel Ridge Regression .
9 %-----%
10
11 close all;
12 clear;
13 clc;
14
15 addpath( '../library' );
16
17 load( 'boston.mat' );
18
19 %% Ex10(c) and (d)
20
21 % Create gamma vector
22 gamma_exp = -40:-26;
23 gamma = 2.^gamma_exp;
24
25 % Create sigma vector
26 sigma_exp = 7:0.5:13;
27 sigma = 2.^sigma_exp;
28
29 % 20 trials
30 for j=1:20
31
32 % Random sampling of data
33 data_train = datasample(boston,round(2/3*size(boston,1)),1,'Replace',false);
34 data_test = setdiff(boston, data_train, 'rows');
35
36 % Training set
37 x_train_j = data_train(:,1:13);
38 y_train_j = data_train(:,14);
39
40 %Test set
41 x_test_j = data_test(:,1:13);
42 y_test_j = data_test(:,14);
43
44 %Dimension of training set x values
45 [dim1, dim2] = size(x_train_j);
```

```

46
47     k_mse = [];
48 % Perform 5-fold cross validation
49 for k=1:5
50
51     k_width = round(dim1/5);
52
53     % Segment x-train data for training
54     x_train = cat( 1, x_train_j(1:(k-1)*k_width , :) , x_train_j(k*k_width
55         +1 :end,:));
56     % Segment x-train data for validation
57     x_valid = x_train_j((k-1)*k_width + 1:k*k_width , :);
58
59     % Segment y-train data for training
60     y_train = cat( 1, y_train_j(1:(k-1)*k_width , :) , y_train_j(k*k_width
61         +1 :end,:));
62     % Segment y-train data for validation
63     y_valid = y_train_j((k-1)*k_width + 1:k*k_width , :);
64
65     [train_dim1 , train_dim2] = size(x_train);
66
67     K = zeros(dim1,dim1,size(sigma,2),1);
68 % From Gram matrix
69 for s=1:size(sigma,2)
70     K(:,:,s) = calK(x_train , x_valid , sigma(s));
71 end
72
73 %Vary selection of gamma
74 for g=1:size(gamma,2)
75     %Vary selection of sigma
76     for ss=1:size(sigma,2)
77         current_k = K(:,:,ss);
78         alpha = kridgereg(current_k(1:train_dim1 ,1:train_dim1 ),y_train
79             ,gamma(g));
80         %Calculate resultant mse for each fold ,each gamma and each
81         %sigma selection and store it into 3D variable
82         k_mse(g,ss,k) = dualcost(current_k(train_dim1+1:end ,1:
83             train_dim1) ,y_valid ,alpha);
84     end
85 end
86 end
87
88 % Calculating the mean MSE over the 5 folds of data
89 mean_k_mse = mean(k_mse ,3);
90
91 %% Ex10 (c)
92 %plot the cross validation error (once as asked in 10_c)
93 if(j == 1)
94     surf(sigma' ,gamma' ,mean_k_mse)
95     colorbar;

```

```

91      set(gcf, 'Color', 'w');
92      xlabel('sigma', 'FontSize', 15);
93      ylabel('gamma', 'FontSize', 15);
94      zlabel('MSE_{cross-validation}', 'FontSize', 15);
95      set(gca, 'YMinorTick', 'on');
96      grid minor
97      axis tight;
98      print('ex10_cv_mse', '-depsc');
99      close all;
100 end
101
102 % Find the combination of gamma and sigma that produces smallest mse
103 [min_gamma_indx, min_sigma_indx] = find(mean_k_mse == min(min(mean_k_mse)));
104
105 % Using optimal gamma and sigma values to calculate optimal alpha
106 % values
107 K_optimal = calK(x_train_j, x_test_j, sigma(min_sigma_indx));
108 alpha_optimal = kridgereg(K_optimal(1:dim1, 1:dim1), y_train_j, gamma(
109     min_gamma_indx));
110
111 %training error
112 mse_training(j) = dualcost(K_optimal(1:dim1, 1:dim1), y_train_j,
113     alpha_optimal);
114
115 %testing error (Using K-test from lower left hand corner of gram matrix)
116 mse_testing(j) = dualcost(K_optimal(dim1+1:end, 1:dim1), y_test_j,
117     alpha_optimal);
118
119 end
120
121 % Calculate the mean and standard deviation of resultant mse of training
122 % and test set
123 avg_mse_training = mean(mse_training)
124 avg_mse_testing = mean(mse_testing)
125
126 std_mse_training = std(mse_training)
127 std_mse_testing = std(mse_testing)

```

dualcost.m

```

1 function mse = dualcost(K,y,a)
2
3 mse = (K*a - y)'*(K*a - y) / size(y,1);

```

kridgereg.m

```

1 function a = kridgereg(K, y, gamma)
2
3 a = (K + (gamma*size(y,1)*eye(size(y,1)))) \ y;

```

Question 4 (Part II): Code**question4_1.m (least squares)**

```

1 %-----%
2 % Module: GI01 – Supervised Learning
3 % Assignment : Coursework 1
4 % Author : Russel Daries , Nitish Mutha
5 % Student ID: 16079408, 15113106
6 % Question: 4 (Least squares)
7 % Section: Part 2
8 % Description: Least Squares sample complexity
9 %-----%
10
11 clear all
12 close all
13
14 % N feaures
15 N = 100;
16 % M samples
17 M = 100;
18 testIterations = 100;
19 testSize = 500;
20
21 avg_missclassified = zeros(testIterations ,N) ;
22 m_error = zeros(N,1) ;
23 %least square
24 for n = 1:N
25     for m = 1:M
26
27         %train
28         X = 2*randi([0 1] , m, n) - 1;
29         y = X(:,1);
30         W = pinv(X)*y;
31
32         %test
33         missClassified = zeros(1,testIterations);
34         for i = 1:testIterations
35             X_test = 2*randi([0 1] , testSize , n) - 1;
36             y_test = X_test(:,1);
37             y_predict = sign(X_test*W);
38             missClassified(1,i) = sum(y_predict~=y_test);
39         end
40         percentError = mean(missClassified*100/testSize);
41         if percentError <= 10
42             m_error(n) = m;
43             avg_missclassified(:,n) = missClassified';
44             break;
45         end
46     end

```

```

47 end
48
49 figure;
50 plot(m_error);
51 set(gcf, 'Color', 'w');
52 xlabel('n features', 'FontSize', 15);
53 ylabel('m samples', 'FontSize', 15);
54 leg=legend('m sample < 10% generalisation error', 'Location', 'Best');
55 set(leg, 'FontSize', 15);
56 set(gca, 'YMinorTick', 'on');
57 grid minor
58 axis tight;
59 print('ex4_least_sq', '-depsc');

```

question4_2.m (perceptron)

```

1 % %
2 % Module: GI01 – Supervised Learning %
3 % Assignment : Coursework 1 %
4 % Author : Russel Daries , Nitish Mutha %
5 % Student ID: 16079408, 15113106 %
6 % Question: 4 (Perceptron) %
7 % Section: Part 2 %
8 % Description: Perceptron sample complexity %
9 % %

10
11 clear all
12 close all
13
14 N = 100;
15 M = 500;
16 testIterations = 100;
17 testSize = 500;
18
19 m_error = zeros(N,1);
20 %perceptron
21 for n = 1:N
22     for m = 1:M
23
24         %train
25         X = 2*randi([0 1], m, n) - 1;
26         y = X(:,1);
27         W = zeros(1,n);
28         for j=1:size(X,1)
29             y_pred = sign(X(j,:)*W);
30             if y_pred*y(j,:) <= 0
31                 W = W + y(j,:)*X(j,:);
32             end
33     end

```

```

34
35 %test
36 missClassified = zeros(1,testIterations);
37 for i = 1:testIterations
38     X_test = 2*randi([0 1], testSize , n) - 1;
39     y_test = X_test(:,1);
40     y_predict = sign(X_test*W);
41     missClassified(1,i) = sum(y_predict~=y_test);
42 end
43 percentError = mean(missClassified)*100/testSize;
44 if percentError <= 10
45     m_error(n) = m;
46     break;
47 end
48 end
49 end
50
51 figure;
52 plot(m_error);
53 set(gcf, 'Color', 'w');
54 xlabel('n features', 'FontSize',15);
55 ylabel('m samples', 'FontSize',15);
56 leg=legend('m sample < 10% generalisation error', 'Location', 'Best');
57 set(leg, 'FontSize',15);
58 set(gca, 'YMinorTick', 'on');
59 grid minor
60 axis tight;
61 print('ex4_perceptron', '-depsc');

```

question4_3.m (winnow)

```

1 %
2 % Module: GI01 – Supervised Learning %
3 % Assignment : Coursework 1
4 % Author : Russel Daries , Nitish Mutha
5 % Student ID: 16079408, 15113106
6 % Question: 4 (Winnow)
7 % Section: Part 2
8 % Description: Winnow sample complexity %
9 %

10
11 clear all %
12 close all %
13
14 N = 100;
15 M = 1000;
16 testIterations = 300;
17 testSize = 500;
18

```

```
19 m_error = zeros(N,1);
20 %winnow
21 for n = 1:N
22     for m = 1:M
23         X = randi([0 1], m, n);
24         y = X(:,1);
25
26     %train
27     W = ones(1,n);
28     for j=1:size(X,1)
29         if X(j,:)*W < n
30             y_pred = 0;
31         else
32             y_pred = 1;
33         end
34         if sum(y_pred~=y(j,:)) > 0
35             W = W .* power(2,(y(j,:)-y_pred)*X(j,:));
36         end
37     end
38
39 %test
40 missClassified = zeros(1,testIterations);
41 for i = 1:testIterations
42     X_test = randi([0 1], testSize, n);
43     y_test = X_test(:,1);
44     y_predict = X_test*W;
45     missClassified(1,i) = sum(y_predict>=n ~= y_test);
46 end
47 percentError = mean(missClassified)*100/testSize;
48 if percentError <= 10
49     m_error(n) = m;
50     break;
51 end
52 end
53 end
54
55 figure;
56 plot(m_error);
57 set(gcf, 'Color', 'w');
58 xlabel('n features', 'FontSize', 15);
59 ylabel('m samples', 'FontSize', 15);
60 leg=legend('m sample < 10% generalisation error', 'Location', 'Best');
61 set(leg, 'FontSize', 15);
62 set(gca, 'YMinorTick', 'on');
63 grid minor
64 axis tight;
65 print('ex4_winnow', '-depsc');
```

question4_4.m (1NN)

```

1 %-----%
2 % Module: GI01 – Supervised Learning
3 % Assignment : Coursework 1
4 % Author : Russel Daries , Nitish Mutha
5 % Student ID: 16079408, 15113106
6 % Question: 4 (1 nearest neighbour)
7 % Section: Part 2
8 % Description: 1 nearest neighbour sample complexity
9 %-----%
10
11 clear all;
12 close all;
13
14 % N features
15 N = 100;
16 % M samples
17 M = 10000;
18 testIterations = 100;
19
20 avg_missclassified = zeros(N,testIterations);
21 m_error = zeros(N,1);
22 %1-NN
23 for n = 1:N
24     for m = 1:M
25         X = 2*randi([0 1], m, n) - 1;
26         y = X(:,1);
27
28         %test
29         missClassified = zeros(1,testIterations);
30         for i = 1:testIterations
31             X_test = 2*randi([0 1], 100, n) - 1;
32             y_test = X_test(:,1);
33             closest_Id = 0;
34             y_predict = zeros(size(X_test,1),1);
35             for j = 1:size(X_test,1)
36                 dist = sqrt(sum((X - X_test(j,:)).^2,2));
37                 [val, closest_Id] = min(dist);
38                 y_predict(j,:) = X(closest_Id,1);
39             end
40
41             missClassified(1,i) = sum(y_predict~=y_test);
42         end
43         percentError = mean(missClassified*100/(100));
44         if percentError <= 10
45             m_error(n) = m;
46             break;
47         end

```

```
48     end
49 end
50
51 figure;
52 plot(m_error);
53 set(gcf, 'Color', 'w');
54 xlabel('n features', 'FontSize', 15);
55 ylabel('m samples', 'FontSize', 15);
56 leg=legend('m sample < 10% generalisation error', 'Location', 'Best');
57 set(leg, 'FontSize', 15);
58 set(gca, 'YMinorTick', 'on');
59 grid minor
60 axis tight;
61 print('ex4_1nn-final', '-depsc');
```

Library Functions

Following are the list of API's used in multiple exercises.

calculateLSR.m

```
1 function lsr = calculateLSR(x,w,n,train_len,test_len)
2 % Function to calculate y_true , weight vector (w) and resultant MSE for
3 % train and test vectors.
4
5 y = x*w + n;
6
7 x_train = x(1:train_len);
8 x_test = x((length(x)-test_len)+1 : end);
9
10 y_train = y(1:train_len);
11 y_test = y(length(x)-test_len+1 : end);
12
13 w_estimator = (x_train'*x_train)\(x_train'*y_train);
14
15 mse_train = meanSquareError(w_estimator,x_train,y_train,train_len);
16 mse_test = meanSquareError(w_estimator,x_test,y_test,test_len);
17
18 lsr = [w_estimator, mse_train, mse_test];
19
20 end
```

calculateLSRex2.m

```
1 function [w_estimator, mse_train, mse_test] = calculateLSRex2(x,w,n,train_len,
2 test_len)
3
4 y = x*w + n;
5
6 % Data segmentation and slicing
7 x_train = x(1:train_len, :);
8 x_test = x((length(x)-test_len)+1 : end, :);
9
10 y_train = y(1:train_len, :);
11 y_test = y(length(x)-test_len+1 : end, :);
12
13 w_estimator = (x_train'*x_train)\(x_train'*y_train);
14
15 % Compute MSE for train and test sets
16 mse_train = meanSquareError(w_estimator,x_train,y_train,train_len);
17 mse_test = meanSquareError(w_estimator,x_test,y_test,test_len);
18
```

19 end

calculateLSRex4.m

```
1 function [w_estimator, mse_train, mse_test] = calculateLSRex4(x,w,n,train_len,
2                                         test_len, gamma)
3 y = x*w + n;
4
5 % Data segmentation and slicing
6 x_train = x(1:train_len, :);
7 x_test = x((length(x)-test_len)+1 : end, :);
8
9 y_train = y(1:train_len, :);
10 y_test = y(length(x)-test_len+1 : end, :);
11
12 % Compute weight vectors
13 w_estimator = (x_train'*x_train + gamma*train_len*eye(size(w,1)))\ (x_train'*
14 y_train);
15
16 % Compute MSE for train and test sets
17 mse_train = meanSquareError(w_estimator,x_train,y_train,train_len);
18 mse_test = meanSquareError(w_estimator,x_test,y_test,test_len);
19
20 end
```

calculateLSRex5.m

```
1 function [w_estimator, mse_train, mse_valid, mse_test] = calculateLSRex5(x,w,n
2 ,train_len, validation_len, test_len, gamma)
3 % Create y vector
4 y = x*w + n;
5
6 % Segment data sets
7 % Training data
8 x_train = x(1:train_len, :);
9 % Validation set
10 x_valid = x(train_len+1 : train_len+1+validation_len, :);
11 % Test Set
12 x_test = x((length(x)-test_len)+1 : end, :);
13
14 % Training set
15 y_train = y(1:train_len, :);
16 % Validation set
17 y_valid = y(train_len+1 : train_len+1+validation_len, :);
18 % Test set
19 y_test = y(length(x)-test_len+1 : end, :);
```

```

19
20 w_estimator = (x_train'*x_train + gamma*train_len*eye(size(w,1)))\ (x_train'*
   y_train);
21
22 % Calculation for MSE of each data set
23 mse_train = meanSquareError(w_estimator,x_train,y_train,train_len);
24 mse_valid = meanSquareError(w_estimator,x_valid,y_valid,validation_len);
25 mse_test = meanSquareError(w_estimator,x_test,y_test,test_len);
26
27
28 end

```

calculateLSRex6.m

```

1 function [w_estimator, mse_train, mse_valid] = calculateLSRex6(x,w,n,train_len
   ,gamma, k)
2
3 y = x*w + n;
4
5 x = x(1:train_len, :);
6 y = y(1:train_len, :);
7
8 [dim1, dim2] = size(x);
9
10 k_width = dim1/5;
11
12 % K-Fold Slicing and segmentation
13 x_valid = x((k-1)*k_width + 1:k*k_width, :);
14 x_train = cat( 1, x(1:(k-1)*k_width, :), x(k*k_width +1 :end,:));
15
16 y_valid = y((k-1)*k_width + 1:k*k_width, :);
17 y_train = cat( 1, y(1:(k-1)*k_width, :), y(k*k_width +1 :end,:));
18
19 % Compute weight vector
20 w_estimator = (x_train'*x_train + gamma*train_len*eye(size(w,1)))\ (x_train'*
   y_train);
21
22 % Compute MSE for train and test sets
23 mse_train = meanSquareError(w_estimator,x_train,y_train,(train_len-k_width));
24 mse_valid = meanSquareError(w_estimator,x_valid,y_valid,k_width);
25
26
27 end

```

calculateLSRTTest.m

```

1 function [w_estimator, mse_test] = calculateLSRTTest(x,w,n,train_len ,
   validation_len ,test_len , gamma)

```

```

2
3 y = x*w + n;
4
5 x_train = x(1:(train_len+validation_len), :);
6 x_valid = x(train_len+1 : train_len+1+validation_len, :);
7 x_test = x((length(x)-test_len)+1 : end, :);
8
9 y_train = y(1:(train_len+validation_len) , :);
10 y_valid = y(train_len+1 : train_len+1+validation_len , :);
11 y_test = y(length(x)-test_len+1 : end, :);
12
13 w_estimator = (x_train'*x_train + gamma*train_len*eye(size(w,1)))\ (x_train'*
    y_train);
14
15 %mse_train = meanSquareError(w_estimator,x_train,y_train,train_len);
16 %mse_valid = meanSquareError(w_estimator,x_valid,y_valid,validation_len);
17 mse_test = meanSquareError(w_estimator,x_test,y_test,test_len);
18
19
20 end

```

meanSquareError.m

```

1 function mse = meanSquareError(w, x, y, len)
2     %Compute MSE for input vectors
3     mse = (1/len) * sum(power((x*w - y),2));
4     %mse = (1/len) * ((w'*x'*x*w) - (2*y'*x*w) + (y'*y));
5 end

```

standard_deviations.m

```

1 function sd = standard_deviations(avg_error,test_errors)
2 %Compute the standard deviation of inputs
3
4 [dim1, dim2] = size(test_errors);
5
6 sum_of_sq = 0;
7 for i=1:dim1
8     sum_of_sq = sum_of_sq + power((test_errors(i) - avg_error),2);
9 end
10 %Output the resultant standard deviation
11 sd = sqrt(sum_of_sq / dim1);

```

calK.m

```

1 function K = calK(x_train,x_test,sigma)
2
3 x = [x_train;x_test];

```

```
4
5 for ii = 1 : size(x,1)
6     for jj = 1 : size(x,1)
7         K(ii,jj) = exp(-(norm(x(ii,:)-x(jj,:),2))^2/(2*sigma^2));
8     end
9 end
```

learnModel.m

```
1 function w_estimator = learnModel(x,y)
2
3 w_estimator = (x'*x)\(x'*y);
4
5 end
```

References

- [1] N. Cristianini and J. Shawe-Taylor, *An Introduction to Support Vector Machines and Other Kernel-based Learning Methods*. Cambridge University Press, 1 ed., 2000.
- [2] MATHWORKS, *MATLAB Matrix Inverse*, 2016 (accessed January 7, 2017).
- [3] U. of Pennsylvania, *CIS520:Machine Learning*, 2016 (accessed January 5, 2017).
- [4] M. Herbster, “Online learning,” October 2016.
- [5] M. Collins, “Convergence proof of perceptron algorithm,” 2010. University of Columbia.
- [6] A. Ng, “The perceptron and large margin classifiers,” 2003. Slides of lecture CSS 229, Stanford University.
- [7] V. N. Vapnik, *Statistical learning theory*. Adaptive and learning systems for signal processing, communications, and control, New York: Wiley, 1998.
- [8] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein, *Introduction to Algorithms*. The MIT Press, 2009.