

WQD7006 WQD7006 MACHINE LEARNING FOR DATA SCIENCE

SEMESTER 1, 2023/2024

FINAL EXAM PART 1 (GROUP 1)

Name	Matric Number
Kar Hong Sam	S2191926

LECTURER'S NAME : DR. MUHAMMAD SHAHREEZA SAFIRUZ BIN KASSIM

DATE OF SUBMISSION : 13 JANUARY 2024

MARKS :

$$|a\rangle h_{\phi}(x) = \theta_{o} + \theta_{1}x_{1} + \theta_{3}x_{2}$$

b) Example (:
$$x_1 = 3$$
, $x_2 = 5$, $[\theta_0 = 1, \theta_1 = 3, \theta_2 = 2]$
 $h_0(y) = 1 + 3(3) + 2(5)$
 $= 20$

Example 2:
$$\chi_1 = 7$$
, $\chi_2 = 2$
 $h_6(x) = 1 + 3(7) + 2(3)$
= 26 m

Example 3:
$$x_1 = 4$$
, $x_2 = 1$

$$h_{\sigma}(x) = 1 + 3(4) + 2(1)$$

() MSE training
$$y = [18, 28]$$
, training $\hat{y} = [20, 26]$
training MSE = $\frac{1}{2}[(18-20)^2+(28-26)^2]$
= $\frac{1}{2}[4+4]$
= $\frac{1}{2}[4+4]$
testing $y = [18]$, toothy $\hat{y} = 15$
testing MSE = $(18-15)^2$
= 9

- d) i) Logistics Regression (LR). Because Logistics Regression is a binary classifier that predict values between I and O. Linear Regression only for numerical values prediction.
 - is Logistis Regression (LR) torm a linear separable surface. Since the data points in Figure above is linear separable honce it can be probably classiful.
 - (ii) No. Since the deta pant is not linearly separable, hence no logistics regression classifler can perfectly classify it.

Question 2:

a. Which classifier has the best generalization performance, i.e., most likely would perform the best when applied to unseen data? Why?

Classifier C, because the testing error rate is the lowest compared to other classifiers indicates well performance in unseen data.

b. Which classifier is underfitting the most? Why?

Classifier A, because it has high training error rate and high testing error compared to other classifiers.

c. Which classifier is overfitting the most? Why?

Classifier B, because it has low training error rate but high testing error rate compared to other classifiers.

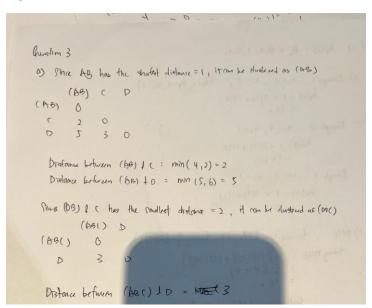
d. Given a dataset, assume we are using linear regression to make a prediction. The dataset is split randomly into training and testing. We increase the training set size gradually. What happens to the mean training and mean testing errors, that is, increase or decrease?

Due to the law of large sample size, the mean for training and testing error will tend to decrease. With more samples to learn the features of data, the model has a better chance to generalize the pattern and identify the relationships in the data. Hence, it will result in decreasing the mean of training and testing errors.

Question 3:

Please this colab link to access the dendogram for both clustering: https://colab.research.google.com/drive/17ROC3SJLITpkvE2YYhU8HSTRJc8rO5JC?usp=sharing

a. Use Single Linkage



b. Use Complete Linkage

```
(bb) (b) (b)

(ba) (b)

(ba) (b)

(ba) (b)

(ba) (c)

(ba) (c)
```