

American University of Sharjah  
School of Engineering  
Department of Computer Engineering  
P. O. Box 26666  
Sharjah, UAE

Instructor: Alex Aklson  
Office: ESB-2172  
Phone: 971-6-515 4893  
E-mail: aaklson@aus.edu  
Semester: Fall 2024

**MLR503 – Data Mining and Knowledge Discovery**  
**Midterm Exam**

**October 31, 2024 (2 hours)**

**Student Name:** \_\_\_\_\_ **Student ID:** \_\_\_\_\_

**Instructions:**

- Write both your *name and ID* above.
- Read the questions carefully; write your answers clearly in the space provided.
- This is a **closed book and closed notes exam**. Only use of calculator is permitted.

**For Instructor's Use Only:**

/ 10	MCQs	Q1
/ 10	Data Types	Q2
/ 20	Linear Regression	Q3
/ 20	Logistic Regression	Q4
/ 15	Decision Trees – Implementation	Q5
/ 5	Decision Tree Ensemble – Theory	Q6
/ 80	<b>Total</b>	

### Question 1 (10 marks)

1. (1 mark) Which of the following is an example of multivariate linear regression?
  - a) Predicting the amount of rainfall based on the temperature.
  - b) Predicting a person's monthly electricity bill based on the number of residents, square footage, and average daily usage.
  - c) Predicting a student's exam score based on their study hours.
  - d) Predicting whether an email is spam or not based on its word count.
  
2. (1 mark) Logistic regression is best suited for which type of problem?
  - a) Predicting continuous outcomes.
  - b) Classifying data into binary categories.
  - c) Finding clusters in unlabeled data.
  - d) Predicting the price of a house based on features such as size, number of floors, and age.
  
3. (1 mark) What happens if the learning rate is too large in gradient descent?
  - a) The algorithm may converge too slowly.
  - b) The cost function may overshoot the minimum, increase or oscillate.
  - c) The algorithm will always find the optimal solution.
  - d) The cost function will not decrease at all.
  
4. (1 mark) Which of the following is true regarding the normal equation and gradient descent? (Select all that apply)
  - a) The normal equation requires iterations, while gradient descent does not.
  - b) Gradient descent is faster for very large datasets compared to the normal equation.
  - c) The normal equation requires solving a matrix inversion.
  - d) Gradient descent always converges faster than the normal equation.

5. **(1 mark)** Which of the following best explains why the linear regression cost function is not suitable for logistic regression?
- a) It is not convex for logistic regression.
  - b) It does not penalize errors for classification.
  - c) It leads to probabilities greater than 1.
  - d) It causes the model to overfit the data.
6. **(1 mark)** Which of the following methods will guarantee finding the global minimum of the cost function for linear regression with gradient descent?
- a) Using a small learning rate.
  - b) Starting with randomly initialized parameters.
  - c) Running gradient descent for an infinite number of iterations.
  - d) The cost function for linear regression is a parabola or a convex function, so gradient descent always finds the global minimum regardless of the learning rate or initialization.
7. **(1 mark)** In logistic regression, which of the following is true about the cost function?
- a) It is non-convex and can have multiple local minima.
  - b) It is convex, ensuring a single global minimum.
  - c) It increases indefinitely with the number of features.
  - d) It can also be minimized using the Normal Equation.
8. **(1 mark)** How does logistic regression handle non-linearly separable data?
- a) By transforming the features using the sigmoid function.
  - b) By increasing the number of iterations in gradient descent.
  - c) By using higher-order polynomial features.
  - d) By reducing the decision threshold.

9. **(1 mark)** How can you extend logistic regression for multi-class classification?
- a) Use a different threshold for each class.
  - b) Use one-vs-all or one-vs-rest approach.
  - c) Train a binary classifier for each pair of classes.
  - d) No extension is necessary, as logistic regression can naturally handle multi-class classification.
10. **(1 mark)** What is the primary purpose of regularization in linear regression?
- a) To improve feature scaling.
  - b) To prevent overfitting by penalizing large coefficients.
  - c) To reduce the dataset size.
  - d) To maximize R-squared value.

## Question 2 (10 marks)

Below is a list of attributes.

Measurement Scale	Type	Attribute
		Number of Children in a Family
		Year of Birth
		Martital Status (e.g., Single, Married, Divorced)
		Purchase Amount (in \$)
		Income Level (e.g., Low Income, Hight Income)

- a) **(5 marks)** For each attribute, identify whether it is categorical (nominal or ordinal) or numerical (continuous or discrete). If the attribute is numerical, also specify the measurement scale (interval or ratio).
- b) **(2 marks)** For the attribute “Income Level”, explain how you would encode it for use in a machine learning model.

- c) **(3 marks)** Describe how you would handle the “Marital Status” attribute if you were to include it in a regression model. Discuss a potential issue and how to address it.

### Question 3 (20 marks)

You are a data scientist working at a real estate company. Your manager asked you to help them price a house for sale that has a lot size of 1500 sq. ft. You ask your manager for a sample data, and you are provided with the following:

Price (\$1000)	Size (sq. ft)	House Number
180	850	House 1
310	1400	House 2
450	2000	House 3

- a) (5 marks) You decide to first normalize your data to have 0 mean and 1 standard deviation. Normalize the data and enter the normalized values in the table below:

Price (\$1000)	Size (sq. ft)	House Number
		House 1
		House 2
		House 3

- b) **(8 marks)** Use Ordinary Least Squares (OLS) to build your model using the **normalized data** and write the equation of the resulting hypothesis. Recall that the OLS equations are the following:

$$w_1 = \frac{\sum_{i=1}^m (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^m (x_i - \bar{x})^2}$$

$$w_0 = \bar{y} - w_1 \bar{x}$$



- c) **(2 marks)** Predict the price that you would recommend to your manager using your model.
- d) **(5 marks)** Your manager is not happy with the price, and wants to sell the house for \$350,000. By leveraging the size coefficient in your hypothesis, what can you recommend the manager do to the house to bring its price to \$350,000 from the price that your model predicted.
- Hint: Consider how the coefficient translates to the original scale of the "Size" feature.*

#### Question 4 (20 marks)

You are a data scientist at a company that provides email security solutions. The company has developed two logistic regression models, *Model A* and *Model B*, to classify incoming emails as “Spam” or “Not Spam”. Both models have been tested on the same dataset, and their performance is summarized in the confusion matrices below.

*Model A:*

		Predicted	
		Not Spam	Spam
Actual	Not Spam	80	20
	Spam	15	85

*Model B:*

		Predicted	
		Not Spam	Spam
Actual	Not Spam	90	10
	Spam	25	75

- a) (5 marks) Calculate the Precision and Recall for both *Model A* and *Model B*. Show your calculations.
- b) (5 marks) Your company is considering deploying one of these models based on different business priorities. For each of the following scenarios, indicate which model (*Model A* or *Model B*) would be more suitable and justify your choice.

Scenario 1:

The company wants to minimize the number of legitimate emails incorrectly marked as spam to avoid inconveniencing users.

Scenario 2:

The company aims to maximize the detection of all spam emails to protect users from potential threats, even if it means some legitimate emails are mistakenly marked as spam.

- c) **(10 marks)** Assume that for your company the cost of a false positive is \$2 per incident, and the cost of a false negative is \$10 per incident. Calculate the total cost associated with each model based on the confusion matrices provided. Which model results in a lower total cost?

### Question 5 (15 marks)

You are a data scientist working for a company that sells home solar panel systems. The marketing team wants to predict whether a homeowner will purchase a solar panel system based on certain characteristics. By understanding customer behavior, the company aims to target its sales efforts more effectively.

You are given a dataset containing information about previous homeowners who were offered solar panel systems, that looks like the following:

Purchased	High Credit Score	Home Ownership	Has Garage	Customer ID
Yes	Yes	Own	Yes	1
Yes	No	Own	Yes	2
No	Yes	Rent	Yes	3
No	No	Rent	Yes	4
Yes	Yes	Own	No	5
No	No	Own	No	6
No	Yes	Rent	No	7
No	No	Rent	No	8
Yes	Yes	Own	Yes	9
Yes	Yes	Own	No	10

- a) (14 marks) Using entropy as the measure of impurity at each node, build a decision tree classifier that can predict whether a customer will purchase a solar panel system. Recall that entropy is calculated as,

$$H(p_1) = - \sum_{i=1}^n p_i \log_2(p_i)$$

where  $p_1$  is the fraction of datapoints belong to class 1.

b) (1 mark) Given a new homeowner with the following characteristics:

- Doesn't have a garage,
- Rents, and,
- Doesn't have a high credit score,

use your decision tree to predict whether this customer will purchase the solar panel system.

### Question 6 (5 marks)

- a) (1 marks) Discuss a serious limitation of a single decision tree.
- b) (4 marks) Explain the difference between bagging and boosting in the context of creating decision tree ensembles. Provide an example of a bagging decision tree ensemble and a boosting decision tree ensemble.