Attempt All Questions

Question 1 : Cluster Analysis

- i. (2 Marks) What is the purpose of a cluster analysis?
- ii. (2 Marks) A discriminant analysis is similar to a cluster analysis; however, there is one fundamental difference. Explain this difference.
- iii. (2 Marks) What is the difference between a linkage method and a distance measure?
- iv. (2 Marks) Compare and contrast any two linkage methods.
- v. (2 Marks) How do we determine the appropriate number of clusters? Give two different visualization methods that are used to display the outcome of a cluster analysis.
- vi. (2 Marks) Compare and contrast agglomerative and divisive clustering.
- vii. (2 Marks) Explain the difference between Ward's method and k means clustering.
- viii. (2 Marks) In the context of hierarchical cluster analysis, distinguish between agglomerative clustering and divisive clustering.
 - ix. (2 Marks) What is a vertical icicle plot used for? Give a brief description, supporting your answer with sketches.
 - x. (2 Marks) Compute the Euclidean distance between the following points.

$$A = (4, 6, 8, 2)$$

$$B = (3, 6, 1, 6)$$

Question 2a: Dimensionality Reduction

The following questions relate to Principal Component Analysis and Factor Analysis

- i. (2 Marks) What is the purpose of a principal component analysis?
- ii. (1 Marks) Principal Component Analysis is a Dimensionality Reduction technique. Explain what this term means.
- iii. (5 Marks) What is meant by the "true" dimension of the data? How does an analyst determine the appropriate number of principal components to retain, making reference to three different approaches.
- iv. (2 Marks) What problems occur if a principal component analysis is carried out on a data matrix where the columns contain measurements on very difference scales? What can be done to overcome this problem?
- v. (2 Marks) The Kaiser-Meyer-Olkin (KMO) statistic is used to measure a certain characteristic of the data. What is this characteristic? Explain how the KMO statistic should be interpreted.
- vi. (2 Marks) Briefly describe the Bartlett Test for Sphericity, with reference to the null and alternative hypotheses, and how those statements relate to the purpose of the test.
- vii. (1 Mark) What is the theoretical difference between principal components analysis and factor analysis?

Question 2b: Multicollinearity

The following questions relate to multicollinearity in the context of multiple regression analysis.

- i. (1 Mark) Define multicollinearity.
- ii. (2 Marks) State two ways in which a multiple regression analysis could be affected by severe multicollinearity.
- iii. (2 Marks) State two ways of formally diagnosing the severity of multicollinearity, making reference to how both should be used to make decisions about the data.

Question 3a: Discriminant Analysis

- i. (2 Marks) What is the purpose of a discriminant analysis?
- ii. (3 Marks) How does discriminant analysis differ from MANOVA?
- iii. (3 Marks) Discuss the condition that determines whether or not a linear or quadratic discriminant rule should be used in a discriminant analysis.
- iv. (3 Marks) Explain the following terms: confusion matrix, prior probabilities cost of misclassification, and apparent error rate.
- v. (1 Mark) The apparent error rate calculated when all observations are used to construct the discriminant rules is known to underestimate the true error rate. What can be done to overcome this problem?

Question 3b: Linear Models

The following questions relate to model selection and validation in the context of multiple regression analysis.

- i. (1 Marks) Explain what variable selections are used for.
- ii. (3 Marks) Compare and contrast the following variable selection procedures
 - a. Forward Selection
 - b. Backward Elimination
 - c. Stepwise Regression
- iii. (1 Mark) Briefly describe how the Akaike Information Criterion would be used in the context of model selection.
- iv. (3 Marks) Describe the process of model validation, with reference to the training, validation and testing phases.

Question 4a: Missing Data

- i. (2 Marks) What is Missing Data? Discuss the implications of Missing Data in the context of a statistical analysis.
- ii. (3 Marks) Compare and contrast the following types of missing data: Missing At Random, Missing Not At Random, Missing Completely at Random.
- iii. (3 Marks) Discuss some of the traditional techniques for dealing with Missing Data, making reference to the limitations of each.

Question 4b: Logistic Regression

- i. (3 Marks) Under what circumstances would you use Logistic Regression?
- ii. (2 Marks) Suppose that, out of a sample of 100 men and 100 women, 70 men drank alcohol in the last week, while 25 women drank alcohol in the past week. Compute the odds ratio for women to men.
- iii. (3 Marks) What is a logit? How can you transform a logit into a probability?
- iv. (2 Marks) What is a dummy variable? Explain how it is used in Logistic Regression. Support your answer with an example.
- v. (2 Marks) Describe how the Likelihood Ratio Test is used for variable selection in Logistic Regression.