



( Semester VI )

- (vii) Choose the correct SQL command to retrieve all columns from a table named 'employees'.
- a) SELECT \* FROM employees;                      b) GET ALL FROM employees;  
c) SHOW employees;                                  d) FETCH \* employees;
- (viii) Identify a common statistical measure used to describe data spread.
- a) Variance    b) Mean  
c) Mode    d) Frequency
- (ix) Identify a key characteristic of NoSQL databases.
- a) They allow flexible data models                      b) They only store relational data  
c) They require predefined schemas                      d) They are slower than SQL databases
- (x) A cricketer's scores in five matches are: 45, 38, 50, 42, 55. What is the range?
- a) 16    b) 17  
c) 22    d) 25
- (xi) Calculate the median of the data: 3, 7, 9, 12, 15, 18, 21
- a) 9    b) 10  
c) 11    d) 12
- (xii) Which of the following is not the correct syntax for creating a set in Python?
- a) set([[1,2],[3,4]])                                      b) set([1,2,2,3,4])  
c) set((1,2,3,4))    d) {1,2,3,4}
- (xiii) Which of the following statements do not align with best practices in data processing and dimensionality reduction?
- a) Training and testing data must be processed differently.                      b) Test transformation is often imperfect.  
c) The primary goal of PCA is statistical analysis, while the secondary goal is data compression.                      d) All of the mentioned
- (xiv) Which of the following statements does not align with standard data science and statistical practices?
- a) ROC curve stands for Receiver Operating Characteristic.                      b) For time series analysis, data must be processed in chunks.  
c) Random sampling must always be done with replacement.                      d) None of the mentioned.
- (xv) Which of the following statements does not correctly describe clustering or classification techniques?
- a) k-means clustering is a method of vector quantization.                      b) k-means clustering aims to partition n observations into k clusters.  
c) k-nearest neighbor is the same as k-means.                      d) k-nearest neighbor is the same as k-means.

### Group-B

(Short Answer Type Questions)

3 x 5=15

2. A dataset follows a normal distribution with a mean of 50 and a standard deviation of 5. Calculate the probability of getting a value greater than 60 using the empirical rule. (3)
3. A company wants to minimize the cost function  $C(x,y)=x^2+y^2$ , subject to the constraint that the total resource allocation satisfies  $x^2 + y^2 = 4$ . Use the Lagrange multiplier method to find the optimal values of x and y. (3)
4. Calculate the Euclidean distance between the two points (2,3) and (5,7) in a k-NN model. (3)
5. Illustrate how to implement multiple linear regression in Python using the scikit-learn library with a sample dataset. (3)
6. A company is using k-means clustering for customer segmentation but is unsure about the optimal number of clusters. Critique whether the Elbow Method or Silhouette Score is a better approach to determine the optimum number of clusters and justify your reasoning. (3)

OR

A company is using Akaike Information Criterion (AIC) and Adjusted R-squared to select the best multiple regression model. Critique which metric is more effective for model selection and justify your reasoning. (3)

**Group-C**  
(Long Answer Type Questions)

5 x 6=30

7. Evaluate the roles of descriptive and inferential statistics in data analysis. Justify the significance of each through practical examples. (5)
8. Draw and explain the bias-variance tradeoff in machine learning. Interpret its impact on model performance. (5)
9. Use Bayes's theorem to solve a probability-based problem and explain the steps involved with a real-life example. (5)
10. Evaluate the performance of different classification models (e.g., Logistic Regression, Random Forest, and SVM) on a given dataset using precision, recall, and F1-score. (5)
11. A regression analysis between apples (y) and oranges (x) resulted in the following least-squares line:  $y = 100 + 2x$ . Predict the implication if oranges are increased by 1 (5)
12. Critically evaluate the different types of biases encountered during the sampling process by analyzing their impact on research validity. Support your evaluation with real-life examples for each type of bias. (5)

**OR**

Critically evaluate the concept of degrees of freedom (DF) in statistics by explaining its significance and impact on the validity of statistical tests. Support your evaluation with a practical example that illustrates its role in hypothesis testing. (5)

\*\*\*\*\*