

KATHMANDU UNIVERSITY
End Semester Examination
February, 2025

Marks Scored:

Level : B.E./B.Sc.
Year : IV

Course : COMP 482
Semester : I

Exam Roll No. :

Time: 30 mins.

F. M. : 10

Registration No.:

Date 04 FEB 2025

SECTION "A"

[20Q. × 0.5 = 10 marks]

Choose the most appropriate answer and **encircle**.

- Which of the following is NOT a data mining task?
 - Entering data into the system manually from printed documents.
 - Detecting and removing noisy data automatically.
 - Identifying new items in a stream of data items.
 - Finding groups of similar items.
- Which of the following visualizations is used for representing distributions?
 - Scatter Plot
 - Line Chart
 - Heatmap
 - Box Plot
- What does numerosity reduction aim to achieve?
 - Reduce the dimensionality of the dataset.
 - Replace the original data with a smaller representation that preserves trends and patterns.
 - Cluster data points into fewer groups.
 - Increase the number of features.
- A confusion matrix shows 90 true positives, 30 false positives, 10 false negatives, and 70 true negatives. What is the accuracy of the model?
 - 0.5
 - 0.2
 - 0.6
 - 0.8
- In k-Nearest Neighbors (k-NN) outlier detection, a data point is considered an outlier if
 - Its distance from the k-nearest neighbors is small.
 - It has many neighboring data points within a small radius.
 - It has few neighboring data points or is far from the neighbors.
 - Its z-score is high.
- In the Apriori algorithm, what is the significance of the minimum support threshold?
 - It determines the confidence level of a rule.
 - It ensures that only frequent itemsets are considered.
 - It decides the number of clusters.
 - It reduces dimensionality.
- Which of the following is a disadvantage of decision trees?
 - Tendency to overfit the data.
 - Lack of interpretability.
 - Inability to handle categorical data.
 - Poor performance on small datasets.

8. In a Bayesian Network, what does each node represent?
- A variable
 - A data point
 - A conditional probability table
 - A dependency
9. A population consists of 3 strata with the following sizes and average values:
- Stratum A: 50 individuals, average age = 30 years
 - Stratum B: 100 individuals, average age = 40 years
 - Stratum C: 150 individuals, average age = 50 years
- If a sample of 60 individuals are drawn using proportional stratified sampling, how many individuals should you sample from each stratum?
- Stratum A: 15, Stratum B: 15, Stratum C: 30.
 - Stratum A: 10, Stratum B: 30, Stratum C: 20
 - Stratum A: 20, Stratum B: 20, Stratum C: 20
 - Stratum A: 10, Stratum B: 20, Stratum C: 30
10. What happens in a 5-fold cross validation?
- The train and test subsets contain respectively 80% and 20% of the data.
 - The train and test subsets contain respectively 95% and 5% of the data.
 - The train and test subsets contain exactly the same number of instances.
 - The train and test subsets contain 70% and 30% of the data.
11. A dataset contains 100 transactions with 40 transactions containing milk, 30 containing bread, and 20 containing milk and bread. What is the confidence of the rule milk \rightarrow bread?
- 20%
 - 50%
 - 66.67%
 - 100%
12. What is the primary goal of data science methodologies?
- To develop complex algorithms.
 - To optimize database performance.
 - To focus only on data visualization.
 - To create a systematic approach for solving data-driven problems.
13. Which performance metric is typically used for evaluating logistic regression models?
- Mean squared error
 - R-squared
 - Confusion Matrix
 - Silhouette score
14. What is the main purpose of backpropagation in neural networks?
- To initialize the weights
 - To calculate the output
 - To determine the activation function
 - To update the weights by minimizing the error
15. Which of the following steps is NOT part of the k-means algorithm?
- Initialize k centroids
 - Use dimensionality reduction to optimize cluster assignments.
 - Update the centroids by taking the mean of assigned points.
 - Assign each data point to the nearest centroid.

16. Which of the following is a key characteristic of ROLAP?
- It stores data in a multidimensional array.
 - It cannot handle large volumes of data.
 - It requires precomputed aggregates for analysis.
 - It relies on relational databases to perform OLAP operations.
17. Consider an OLAP cube with three dimensions: Product (P1, P2), Region (North, South), and Time (2023, 2024). The cube stores sales revenue for each combination. If the cube contains $2 \times 2 \times 2 = 8$ data cells, what would happen to the number of cells if a new dimension Customer (C1, C2, C3) is added?
- 12 cells
 - 16 cells
 - 24 cells
 - 32 cells
18. What is the main difference between classification and regression in data mining?
- Classification predicts continuous values, while regression predicts categories.
 - Classification predicts categories, while regression predicts continuous values.
 - Classification groups data into clusters, while regression finds patterns in data.
 - Classification is unsupervised, while regression is supervised.
19. Suppose you are given the following data points in 1D space: {2, 4, 10, 12, 18, 20, 25} and you want to divide them into $k = 2$ clusters. Initially, the centroids are $C1 = 4$, and $C2 = 1$. After one iteration of the k-means algorithm, which of the following will be the new position of the centroids?
- $C1 = 7, C2 = 21$
 - $C1 = 5.33, C2 = 21.67$
 - $C1 = 5.33, C2 = 18.67$
 - $C1 = 6, C2 = 19$
20. What is the key assumption in the DBSCAN (Density-Based Spatial Clustering of Applications with Noise) method for outlier detection?
- Data points are uniformly distributed.
 - Clusters are circular in shape.
 - Data points are normally distributed.
 - Noise points are sparse in regions of low density.

KATHMANDU UNIVERSITY
End Semester Examination
February, 2025

Level : B.E./B.Sc.
Year : IV
Time : 2 hrs. 30mins.

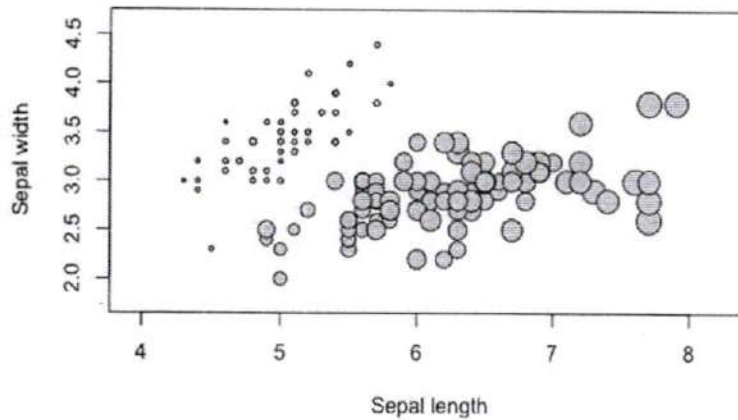
04 FEB 2025

Course : COMP 482
Semester : I
F. M. : 40

SECTION "B"
[6Q. × 4 = 24 marks]

Attempt *ANY SIX* questions.

1. Discuss the Data Science Hierarchy of Needs. Why is it important to follow data mining methodologies? Explain, in brief, any one data mining methodology. [1 + 1 + 2]
2. Discuss the importance of exploratory data analysis. Suppose you are given the iris dataset with 3 attributes - sepal length, sepal width, and petal length. What can you say



- about the dataset from the following plot? [2 + 2]
3. Define dimensions, and measures. Give an example of a constellation schema with at least 2 measures, and 3 dimensions. [2 + 2]
 4. Explain the role of lift and correlation analysis in measuring the interestingness of association rules. Give an example of a situation where lift can better measure rule interestingness than confidence. [2 + 2]
 5. Compare and contrast classification and regression tasks. Design a naive Bayes classifier from Table 1. [2 + 2]
 6. What kind of data can be used with a single-layer neural network for binary classification? Design a perceptron for implementing OR gate. [1 + 3]
 7. Differentiate between (*ANY TWO*): [2 + 2]
 - a. Numerosity reduction and dimensionality reduction
 - b. Divisive and agglomerative clustering
 - c. SEMMA and CRISP-DM

P.T.O.

SECTION "C"
[2Q. × 8 = 16 marks]

Attempt *ANY TWO* questions.

8. Consider the following dataset. Simulate a 3-fold cross validation for evaluating two classification algorithms, clearly explaining how a 3-fold cross-validation works. How would you choose the value of k for k-fold cross-validation? Explain the scenario when you would choose $k=n$, where n is the number of instances in the dataset. [5 + 2 + 1]

Table 1 Dataset for Question 5 and 8

Instance	Feature1	Feature2	Class
1	S	A	+
2	S	B	-
3	S	A	+
4	M	A	+
5	S	B	+
6	M	A	-

9. Construct a dendrogram that shows the hierarchical relationship between the following data points: (3, 1, 1), (2, 2, 2), (1, 2, 5), (3, 2, 4), (4, 5, 5). Use the single link method to compute distance between clusters. Divide the data points into two clusters using the constructed dendrogram, and compute Silhouette coefficient of any one data point. [4 + 1 + 3]
10. What is the main idea behind statistical methods of outlier detection? In the Interquartile range (IQR) method of detecting outliers, a data point is an outlier if it is more than 1.5 times IQR above the third quartile or below the first quartile. Note that IQR is the difference between the third quartile (Q3) and the first quartile (Q1). Using this method, determine if the highest and the lowest scores are outliers in the set of 19 test scores: 5, 7, 7, 8, 9, 10, 11, 11, 11, 11, 13, 14, 14, 14, 16, 16, 16, 17, 20. How would you use this method to find outliers in a multivariate data? [2 + 4 + 2]