



M.TECH DEGREE EXAMINATIONS: DEC 2022

(Regulation 2018)

First Semester

DATA SCIENCE

P18ITI1203: Data Mining Techniques

COURSE OUTCOMES

- CO1:** Explain the techniques for data pre processing
- CO2:** Apply association rules algorithm for correlation analysis
- CO3:** Apply decision tree algorithm for classification
- CO4:** Apply Bayesian networks algorithm for classification
- CO5:** Apply various clustering algorithms for different datasets
- CO6:** Estimate the classifier accuracy with training, testing and cross validation datasets

Time: Three Hours

Maximum Marks: 100

Answer all the Questions:-

PART A (10 x 1 = 10 Marks)

1. The following items consists of two statements, one labeled as the “Assertion (A)” and the other as “Reason (R). You are to examine those two statements carefully and select the answers to these items using the codes given below Codes: CO1 [K₁]

Assertion (A): A relational database is a collection of tables, each of which is assigned a unique name.

Reason (R): Each table consists of a set of attributes (*columns* or *fields*) and usually stores a large set of tuples (*records* or *rows*).

- | | |
|---|---|
| a) both A and R are individually true and R is the correct explanation of A | b) both A and R are individually true but R is not the correct explanation of A |
| c) A is true but R is false | d) A is false but R is true. |
2. A _____ is a set of mathematical functions that describe the behavior of the objects in a target class in terms of random variables and their associated probability distributions. CO1 [K₁]
- | | |
|-----------------------|----------------------|
| a) Hierarchical model | b) Descriptive model |
| b) Statistical model | d) Linear model |
3. _____ methods smooth a sorted data value by consulting its “neighborhood,” that is, the values around it. CO2 [K₁]
- | | |
|---------------|---------------------|
| a) Regression | b) Outlier analysis |
| c) Extraction | d) Binning |

4. Matching type item with multiple choice code

CO4 [K₁]

List I	List II
A. Nominal attribute	i. Boolean
B. Binary Attributes	ii. Integer or real values
C. Ordinal Attributes	iii. Symbols or names of things
D. Numeric Attributes	iv. Meaningful order or ranking

- | | A | B | C | D |
|----|-----|-----|----|-----|
| a) | ii | iii | i | iv |
| b) | iii | i | iv | ii |
| c) | ii | iii | i | iv |
| d) | iv | i | ii | iii |

5. The following items consists of two statements, one labeled as the “Assertion (A)” and the other as “Reason (R). You are to examine those two statements carefully and select the answers to these items using the codes given below Codes:

CO4 [K₁]

Assertion (A): Normalization, where the attribute data are scaled so as to fall within a smaller range, such as 1.0 to 1.0, or 0.0 to 1.0.

Reason (R): Discretization, where the raw values of a numeric attribute (e.g., *age*) are replaced by interval labels (e.g., 0–10, 11–20, etc.) or conceptual labels (e.g., *youth, adult, senior*).

- | | |
|---|---|
| a) both A and R are individually true and R is the correct explanation of A | b) both A and R are individually true but R is not the correct explanation of A |
| c) A is true but R is false | d) A is false but R is true. |

6. _____ techniques can be applied to obtain a reduced representation of the data set that is much smaller in volume, yet closely maintains the integrity of the original data.

CO3 [K₁]

- | | |
|-------------------|------------------------|
| a) Data Reduction | b) Data extraction |
| c) Data Loading | d) Data Transformation |

7. Which of the following are extended data types?

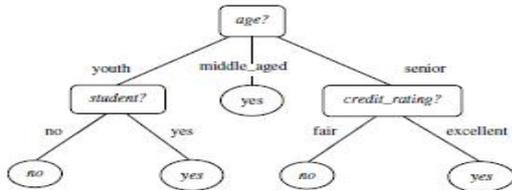
CO2 [K₁]

- | | |
|--|--|
| i) sequential and time-series patterns | ii) structural (e.g., tree, lattice, graph) patterns |
| iii) spatial (e.g., colocation) patterns | iv) approximate patterns |
| a) i, ii, iii | b) i, iii, iv |

- | | |
|--------------|-----------------|
| c) i, ii, iv | d) ii, iii, iii |
|--------------|-----------------|

- (d) Can you find (roughly) the first quartile (Q1) and the third quartile (Q3) of the data?
- (e) Give the five-number summary of the data.
- (f) Show a boxplot of the data.
- (g) How is a quantile–quantile plot different from a quantile plot?

- 23. Differentiate between snowflake schema and fact constellation query model with neat diagram. CO2 [K₂]
- 24. Explain about FP-growth algorithm with an example. CO2 [K₂]
- 25. Explain in detail about metrics for evaluating classifier performance CO3 [K₂]
- 26. How to extract rule from a decision tree? Explain for the below decision tree. CO3 [K₃]



- 27. Describe about k-Means: A Centroid-Based Technique for clustering. CO4 [K₂]
- 28. Differentiate agglomerative versus divisive hierarchical clustering. CO4 [K₂]
- 29. How supervised methods are used for outlier detection methods? CO5 [K₂]
- 30. Explain about proximity-based methods with an example. CO5 [K₂]

**Answer any TWO Questions
PART D (2 x 10 = 20 Marks)**

- 31. Briefly describe about attribute selection measures for the following class labeled training tuples from the *AllElectronics* Customer Database. CO2 [K₃]

RID	age	income	student	credit_rating	Class: buys_computer
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle_aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle_aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle_aged	medium	no	excellent	yes
13	middle_aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

- 32. Enumerate how support vector machines (SVMs) is used for the classification of both linear and nonlinear data. CO3 [K₃]
- 33. Explain in detail about Partitioning methods and Hierarchical methods for clustering. CO5 [K₃]
