

Code: 20IT4501E

III B.Tech - I Semester – Regular Examinations - DECEMBER 2022

**DATA MINING
(INFORMATION TECHNOLOGY)**

Duration: 3 hours

Max. Marks: 70

Note: 1. This paper contains questions from 5 units of Syllabus. Each unit carries 14 marks and have an internal choice of Questions.

2. All parts of Question must be answered in one place.

BL – Blooms Level

CO – Course Outcome

			BL	CO	Max. Marks
UNIT-I					
1	a)	Write a brief note on relational databases and data warehouses.	L2	CO1	5 M
	b)	Describe the data mining functionalities, and the kinds of patterns they can discover.	L2	CO1	9 M
OR					
2	a)	Describe the various phases in knowledge discovery process with a neat diagram.	L2	CO1	10 M
	b)	Explain how the evolution of database technology led to data mining.	L2	CO1	4 M
UNIT-II					
3	a)	What is the curse of dimensionality? How to reduce it?	L3	CO2	7 M
	b)	Illustrate binning methods for data smoothing.	L3	CO2	7 M
OR					
4	a)	In real-world data, tuples with missing values for some attributes are a common occurrence. Describe various methods for handling this problem.	L3	CO2	8 M

	b)	Discuss issues to consider during data integration.	L2	CO2	6 M
--	----	---	----	-----	-----

UNIT-III

5	a)	Can we design a method that mines the complete set of frequent item sets without candidate generation? Explain with example.	L3	CO3	9 M
	b)	Prove that all nonempty subsets of a frequent itemset must also be frequent.	L4	CO3	5 M

OR

6	a)	How are association rules generated from frequent itemsets? Illustrate.	L3	CO3	7 M
	b)	Apply FP-Growth algorithm to the following transactional data to find frequent itemsets. List all frequent itemsets with their minimum support count of 2 and confidence = 50%.	L4	CO5	7 M

TID	List of Item IDs
1	i1,i3,i5,i7
2	i2,i4,i6,i8
3	i1,i3,i5,i7
4	i9,i7,i5,i1
5	i2,i4,i6,i7
6	i1,i2,i3,i4
7	i3,i4,i5,i6
8	i7,i8,i6,i1
9	i8,i5,i3,i2
10	i1,i3,i4,i6

UNIT-IV

7	a)	Define information gain and explain its importance in decision tree induction.	L2	CO4	4 M
---	----	--	----	-----	-----

	<p>b) The following table consists of training data from an employee database. The data have been generalized. For example, “31 : : : 35” for age represents the age range of 31 to 35. For a given row entry, count represents the number of data tuples having the values for department, status, age, and salary given in that row.</p> <table border="1" data-bbox="311 645 1045 1209"> <thead> <tr> <th><i>department</i></th> <th><i>status</i></th> <th><i>age</i></th> <th><i>salary</i></th> <th><i>count</i></th> </tr> </thead> <tbody> <tr><td>sales</td><td>senior</td><td>31...35</td><td>46K...50K</td><td>30</td></tr> <tr><td>sales</td><td>junior</td><td>26...30</td><td>26K...30K</td><td>40</td></tr> <tr><td>sales</td><td>junior</td><td>31...35</td><td>31K...35K</td><td>40</td></tr> <tr><td>systems</td><td>junior</td><td>21...25</td><td>46K...50K</td><td>20</td></tr> <tr><td>systems</td><td>senior</td><td>31...35</td><td>66K...70K</td><td>5</td></tr> <tr><td>systems</td><td>junior</td><td>26...30</td><td>46K...50K</td><td>3</td></tr> <tr><td>systems</td><td>senior</td><td>41...45</td><td>66K...70K</td><td>3</td></tr> <tr><td>marketing</td><td>senior</td><td>36...40</td><td>46K...50K</td><td>10</td></tr> <tr><td>marketing</td><td>junior</td><td>31...35</td><td>41K...45K</td><td>4</td></tr> <tr><td>secretary</td><td>senior</td><td>46...50</td><td>36K...40K</td><td>4</td></tr> <tr><td>secretary</td><td>junior</td><td>26...30</td><td>26K...30K</td><td>6</td></tr> </tbody> </table> <p>(i) How would you modify the basic decision tree algorithm to take into consideration the count of each generalized data tuple (i.e., of each row entry)?</p> <p>(ii) Use your algorithm to construct a decision tree from the given data.</p> <p>(iii) Given a data tuple having the values “systems,” “26 . . . 30,” and “46–50K” for the attributes department, age, and salary, respectively, what would a naive Bayesian classification of the status for the tuple be?</p>	<i>department</i>	<i>status</i>	<i>age</i>	<i>salary</i>	<i>count</i>	sales	senior	31...35	46K...50K	30	sales	junior	26...30	26K...30K	40	sales	junior	31...35	31K...35K	40	systems	junior	21...25	46K...50K	20	systems	senior	31...35	66K...70K	5	systems	junior	26...30	46K...50K	3	systems	senior	41...45	66K...70K	3	marketing	senior	36...40	46K...50K	10	marketing	junior	31...35	41K...45K	4	secretary	senior	46...50	36K...40K	4	secretary	junior	26...30	26K...30K	6	L4	CO4	10 M
<i>department</i>	<i>status</i>	<i>age</i>	<i>salary</i>	<i>count</i>																																																												
sales	senior	31...35	46K...50K	30																																																												
sales	junior	26...30	26K...30K	40																																																												
sales	junior	31...35	31K...35K	40																																																												
systems	junior	21...25	46K...50K	20																																																												
systems	senior	31...35	66K...70K	5																																																												
systems	junior	26...30	46K...50K	3																																																												
systems	senior	41...45	66K...70K	3																																																												
marketing	senior	36...40	46K...50K	10																																																												
marketing	junior	31...35	41K...45K	4																																																												
secretary	senior	46...50	36K...40K	4																																																												
secretary	junior	26...30	26K...30K	6																																																												
OR																																																																
8	a) Consider a school with a total population of 100 persons. These 100 persons can be seen either as ‘Students’ and ‘Teachers’ or as a	L4	CO5	7 M																																																												

		<p>population of ‘Males’ and ‘Females’.</p> <p>With below tabulation of the 100 people, what is the conditional probability that a certain member of the school is a ‘Teacher’ given that he is a ‘Man’?</p> <table border="1" style="margin-left: auto; margin-right: auto;"> <thead> <tr> <th></th> <th>Female</th> <th>Male</th> <th>Total</th> </tr> </thead> <tbody> <tr> <td>Teacher</td> <td>8</td> <td>12</td> <td>20</td> </tr> <tr> <td>Student</td> <td>32</td> <td>48</td> <td>80</td> </tr> <tr> <td>Total</td> <td>40</td> <td>60</td> <td>100</td> </tr> </tbody> </table>		Female	Male	Total	Teacher	8	12	20	Student	32	48	80	Total	40	60	100			
	Female	Male	Total																		
Teacher	8	12	20																		
Student	32	48	80																		
Total	40	60	100																		
	b)	Why is tree pruning useful in decision tree induction? What is a drawback of using a separate set of tuples to evaluate pruning?	L2	CO4	7 M																
UNIT-V																					
9	a)	Compare k-means with k-medoids algorithms for clustering.	L3	CO4	6 M																
	b)	Discuss various evaluation measures used to evaluate clustering algorithms.	L2	CO4	8 M																
OR																					
10	a)	Discuss the similarity measures and distance measures frequently used in clustering the data.	L3	CO4	10 M																
	b)	Discuss about key issues in Hierarchical clustering.	L3	CO4	4 M																