



KIIT Deemed to be University
Online End Semester Examination(Spring Semester-2021)

Subject Name & Code: **Data Mining and Data Warehousing (IT-3031)**
Applicable to Courses:

Full Marks=50

Time:2 Hours

SECTION-A(Answer All Questions. Each question carries 2 Marks)

Time:30 Minutes

(7×2=14 Marks)

<u>Question No</u>	<u>Question Type(MCQ/SAT)</u>	<u>Question</u>	<u>CO Mapping</u>	<u>Answer Key (For MCQ Questions only)</u>
<u>Q.No:1</u>	SAT	What is descriptive and predictive data mining?	CO1	
	SAT	What are the factors that leads to mining the data?	CO1	
	SAT	Data mining is applicable for any kind of information repository. Justify.	CO1	
	SAT	Give examples of incomplete and inconsistent data.	CO1	
<u>Q.No: 2</u>	SAT	Given two objects represented by the tuples (23, 11, 42, 10) and (20, N, 36, 7). Compute the Manhattan distance between the two objects. N = right most significant digit of your Roll No. (Ex:- for Roll No. 180655, N=5)	CO2	
	SAT	Given two objects represented by the tuples (23, 11, 42, 10) and (20, N, 36, 7). Compute the Euclidean distance between the two objects. N = right most significant digit of your Roll No. (CO2	

		Ex:- for Roll No. 180656, N=6)		
	SAT	Suppose that the data for analysis includes the attribute age. The age values for the data tuples are 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 26, 27, 30, N What is the mean of the data? What is the median? N = 20 + right most significant digit of your Roll No. (Ex:- for Roll No. 180652, N=20+2=22)	CO2	
	SAT	Obtain trimmed 20% mean for the following data: 40,36,48,50,52,52,56,60,63, 80.	CO2	
Q.No: 3	SAT	Explain the importance of cross validation.	CO3	
	SAT	How FP growth tree is better than Apriori?	CO3	
	SAT	Does SVM algorithm have high classification accuracy in high dimensional space? Justify your answer.	CO3	
	SAT	What is over fitting and what can you do to prevent it?	CO3	
Q.No: 4	SAT	State the advantages of tree pruning.	CO3	
	SAT	What is meant by spatial database? Mention its features.	CO5	
	SAT	List out two major strengths of decision tree method.	CO2	
	SAT	Data marts are becoming popular day by day. Comment.	CO3	
Q.No: 5	SAT	Differentiate between median and midrange of a data set.	CO4	
	SAT	List the disadvantages of the k-NN classification technique?	CO4	
	SAT	How to improve the performance of the Apriori method for association mining?	CO4	
	SAT	List out some social impacts of Data mining.	CO4	

Q.No: 6	SAT	Differentiate between supervised and unsupervised learning.	CO5	
	SAT	Differentiate between metadata and data mart.	CO5	
	SAT	What is the difference between discrimination and classification?	CO5	
	SAT	List the short comings of K-means algorithm.	CO5	
Q.No: 7	SAT	Explain web mining and write some applications of it.	CO6	
	SAT	Explain spatial mining and write some applications of it.	CO6	
	SAT	Explain text mining and write some applications of it.	CO6	
	SAT	Explain multimedia mining and write some applications of it.	CO6	

SECTION-B(Answer Any Three Questions. Each Question carries 12 Marks)

Time: 1 Hour and 30 Minutes

(3×12=36 Marks)

<u>Question No</u>	<u>Question</u>	<u>CO Mapping</u> <u>(Each question should be from the same CO(s))</u>
Q.No: 8	What is the need of preprocessing the data before mining it? Explain the various techniques use for data transformation and discretization.	CO2
	What are the different schemas to represent multidimensional data models? What is a data cube? Explain with an example. Describe	

	the OLAP operations that can be performed on a data cube.																					
	Explain different types of attributes with examples. What are the different measures to represent dispersion of data? Discuss the purpose of a boxplot.																					
Q.No: 9	<p>What is frequent pattern analysis? Explain the apriori property and the general steps of the Apriori algorithm for association rule mining.</p> <p>The database below has four transactions. Find all frequent itemsets and association rules, if the minimum support is 60% and the minimum confidence is 80%.</p> <p>Trans_id Itemlist</p> <p>T1 {K, A, D, B}</p> <p>T2 {D, A C, E, B}</p> <p>T3 {C, A, B, E}</p> <p>T4 {B, A, D}</p> <p>You are given the transaction data shown in the table below from a fast food restaurant. There are 9 distinct transactions and a total of 5 meal items that are involved in the transactions.</p> <table><tr><th>Order</th><th>List of Item IDs</th></tr><tr><td>1</td><td>M1, M2, M5</td></tr><tr><td>2</td><td>M2, M4</td></tr><tr><td>3</td><td>M2, M3</td></tr><tr><td>4</td><td>M1, M2, M4</td></tr><tr><td>5</td><td>M1, M3</td></tr><tr><td>6</td><td>M2, M3</td></tr><tr><td>7</td><td>M1, M3</td></tr><tr><td>8</td><td>M1, M2, M3, M5</td></tr><tr><td>9</td><td>M1, M2, M3</td></tr></table> <p>Assuming the minimum support is 2/9 (.222) and the minimum confidence is 7/9 (.777),</p> <p>i. Apply the Apriori algorithm to the transactions and identify all frequent itemsets</p>	Order	List of Item IDs	1	M1, M2, M5	2	M2, M4	3	M2, M3	4	M1, M2, M4	5	M1, M3	6	M2, M3	7	M1, M3	8	M1, M2, M3, M5	9	M1, M2, M3	<p>CO4</p> <p>CO4</p>
Order	List of Item IDs																					
1	M1, M2, M5																					
2	M2, M4																					
3	M2, M3																					
4	M1, M2, M4																					
5	M1, M3																					
6	M2, M3																					
7	M1, M3																					
8	M1, M2, M3, M5																					
9	M1, M2, M3																					

	<p>the k-means algorithm to show</p> <p>(a) The three cluster centers after the first round of execution.</p> <p>(b) The final three clusters.</p>	CO3																																																							
Q.No: 11	<p>Suppose that a data warehouse for Big_University consists of the four dimensions: Student, Course, Semester, and Instructor, and two measures count and avg_grade. At the lowest conceptual level (e.g., for a given student, course, semester, and instructor combination), the avg_grade measure stores the actual course grade of the student. At higher conceptual levels, avg_grade stores the average grade for the given combination.</p> <p>(a) Draw a snowflake schema diagram for the data warehouse.</p> <p>(b) Starting with the base cuboid (Student, Course, Semester, Instructor), what specific OLAP operations (e.g., roll-up from semester to year) should you perform in order to list the average grade of CS courses for each Big_University student.</p> <p>(c) If each dimension has five levels (including all), such as “Student < Major < Status < University < all”, how many cuboids will this cube contain (including the base)?</p> <p>Compare and contrast the characteristics of the four major clustering methods. Give examples of each category.</p> <p>Discuss the two types of hierarchical clustering methods. In your opinion, which method is more preferable and why? Explain a dendrogram and its use in representing hierarchical clustering.</p> <p>State the Bayes theorem and explain the working of a Naive bayes classifier? What is the major assumption that the Naive Bayes classifier follows?</p> <table><tr><th>T_id</th><th>Refund</th><th>Marital Status</th><th>Taxable Income</th><th>Evade</th></tr><tr><td>1</td><td>Yes</td><td>Single</td><td>125K</td><td>No</td></tr><tr><td>2</td><td>No</td><td>Married</td><td>100K</td><td>No</td></tr><tr><td>3</td><td>No</td><td>Single</td><td>70K</td><td>No</td></tr><tr><td>4</td><td>Yes</td><td>Married</td><td>120K</td><td>No</td></tr><tr><td>5</td><td>No</td><td>Divorced</td><td>95K</td><td>Yes</td></tr><tr><td>6</td><td>No</td><td>Married</td><td>60K</td><td>No</td></tr><tr><td>7</td><td>Yes</td><td>Divorced</td><td>220K</td><td>No</td></tr><tr><td>8</td><td>No</td><td>Single</td><td>85K</td><td>Yes</td></tr><tr><td>9</td><td>No</td><td>Married</td><td>75K</td><td>No</td></tr><tr><td>10</td><td>No</td><td>Single</td><td>90K</td><td>Yes</td></tr></table> <p>Given the above dataset for tax evaders, use a Naive bayes classifier to predict whether the following customer will evade tax.</p> <p>Y = (Refund = No, Marital Status=Married, Income = 120K)</p>	T_id	Refund	Marital Status	Taxable Income	Evade	1	Yes	Single	125K	No	2	No	Married	100K	No	3	No	Single	70K	No	4	Yes	Married	120K	No	5	No	Divorced	95K	Yes	6	No	Married	60K	No	7	Yes	Divorced	220K	No	8	No	Single	85K	Yes	9	No	Married	75K	No	10	No	Single	90K	Yes	CO5
T_id	Refund	Marital Status	Taxable Income	Evade																																																					
1	Yes	Single	125K	No																																																					
2	No	Married	100K	No																																																					
3	No	Single	70K	No																																																					
4	Yes	Married	120K	No																																																					
5	No	Divorced	95K	Yes																																																					
6	No	Married	60K	No																																																					
7	Yes	Divorced	220K	No																																																					
8	No	Single	85K	Yes																																																					
9	No	Married	75K	No																																																					
10	No	Single	90K	Yes																																																					
		CO4																																																							
		CO5																																																							

