



KIIT Deemed to be University
Online End Semester Examination(Autumn Semester-2020)

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>	<u>MCQ</u>	Given two objects represented by the tuples (22, 1, 42, 10) and (20, 0, 36, 8). Compute the Euclidean distance between the two objects. a) 6.7082 b) 6.6209 c) 5.6034 d) 5.7062	CO1	A
	<u>MCQ</u>	What is the interquartile range for an even sample size (81,62, 77,63,72,64,70,76,81,64) a) Interquartile range (79-64)=15 b) Interquartile range (77-64)=13 c) Interquartile range (76-64)=12 d) Interquartile range (81-64)=17	CO2	B
	<u>MCQ</u>	Which one is correct normalized data after performing min-max normalization by setting min = 0 and max = 1 for the data 200, 300, 400, 600, 1000. a) 0, 0.5, 0.25, 0.125, 1 b) 0, 0.125, 0.25, 0.5, 1 c) 1, 0.5, 0.125, 0.25, 0 d) 1, 0.125, 0.25, 0.5, 0	CO1	B
	<u>MCQ</u>	Suppose that the data for analysis includes the attribute age. 13, 15, 16, 46, 19, 20, 20, 21, 22, 22, 25, 25, 25, 52, 30, 33, 35, 35, 33, 35, 35, 36,40, 45, 16, 25, 70. What is the mode, first quartile and third quartile of the data? a) mode=25, first quartile=22, third quartile=35 b) mode=22,25, first quartile=22, third quartile=35	CO1	D

		c) mode=25, first quartile=20, third quartile=35 d) mode=22,25, first quartile=20, third quartile=35																						
<u>Q.No:2</u>	<u>SAT</u>	Define association and co-relation.	CO3																					
	<u>SAT</u>	What is Jaccard co-efficient? Explain with example.	CO1																					
	<u>SAT</u>	List out the functionality of metadata.	CO2																					
	<u>SAT</u>	What is correlation analysis.	CO3																					
<u>Q.No:3</u>	<u>SAT</u>	What is the need of Information Gain in decision tree classifier?	CO4																					
	<u>SAT</u>	What is discrete and continuous data in data mining?	CO1																					
	<u>SAT</u>	How to handle tuples with missing values for some attributes?	CO1																					
	<u>SAT</u>	Define anti-monotone property.	CO3																					
<u>Q.No:4</u>	<u>MCQ</u>	For d items, there are _____ possible candidate item set and _____ rules can be created. a) pow(3,d), pow(2,d)-pow(3,d+1)+1 b) pow(2,d), pow(3,d)-pow(2,d+1)+1 c) pow(2,d), pow(2,d)-pow(3,d+1)+1 d) pow(3,d), pow(3,d)-pow(2,d+1)+1	CO3	B																				
	<u>MCQ</u>	Which one is incorrect option for support and confidence value for the following transaction data ? <table border="1"><thead><tr><th>TID</th><th>ITEMS</th></tr></thead><tbody><tr><td>1</td><td>Bread, milk</td></tr><tr><td>2</td><td>Bread, Diaper, Beer, Eggs</td></tr><tr><td>3</td><td>Milk, diaper, beer, coke</td></tr><tr><td>4</td><td>Bread, milk, diaper, beer</td></tr><tr><td>5</td><td>Bread, milk, diaper, coke</td></tr></tbody></table> a) {Diaper,Beer} → {Milk} (s=0.4, c=0.67) b) {Milk,Diaper} → {Beer} (s=0.4, c=0.67) c) {Milk} → {Diaper,Beer} (s=0.4, c=0.5) d) {Milk,Beer} → {Diaper} (s=0.4, c=0.6)	TID	ITEMS	1	Bread, milk	2	Bread, Diaper, Beer, Eggs	3	Milk, diaper, beer, coke	4	Bread, milk, diaper, beer	5	Bread, milk, diaper, coke	CO3	D								
TID	ITEMS																							
1	Bread, milk																							
2	Bread, Diaper, Beer, Eggs																							
3	Milk, diaper, beer, coke																							
4	Bread, milk, diaper, beer																							
5	Bread, milk, diaper, coke																							
	<u>MCQ</u>	Perform KNN for "K=3" on the following dataset and generate the class level for the input (Acid durability =3 , strength=7, class=?). <table border="1"><thead><tr><th>Name</th><th>Acid durability</th><th>strength</th><th>class</th></tr></thead><tbody><tr><td>Type 1</td><td>7</td><td>7</td><td>bad</td></tr><tr><td>Type 2</td><td>7</td><td>4</td><td>bad</td></tr><tr><td>Type 3</td><td>3</td><td>4</td><td>good</td></tr><tr><td>Type 4</td><td>1</td><td>4</td><td>good</td></tr></tbody></table> a) Good b) Bad c) Invalid	Name	Acid durability	strength	class	Type 1	7	7	bad	Type 2	7	4	bad	Type 3	3	4	good	Type 4	1	4	good	CO4	A
Name	Acid durability	strength	class																					
Type 1	7	7	bad																					
Type 2	7	4	bad																					
Type 3	3	4	good																					
Type 4	1	4	good																					

		a) none		
	<u>MCQ</u>	Bayesian classifiers is a) A class of learning algorithm that tries to find an optimum classification of a set of examples using the probabilistic theory. b) Any mechanism employed by a learning system to constrain the search space of a hypothesis c) An approach to the design of learning algorithms that is inspired by the fact that when people encounter new situations, they often explain them by reference to familiar experiences, adapting the explanations to fit the new situation. d) None	CO4	A
<u>Q.No:5</u>	<u>SAT</u>	What the condition two item sets A and B have no co-relation between them?	CO3	
	<u>SAT</u>	What is tree pruning explain with example.	CO4	
	<u>SAT</u>	Explain ETL with respect to data warehouse.	CO2	
	<u>SAT</u>	What is Bassel's correction?	CO3	
<u>Q.No:6</u>	<u>MCQ</u>	Knowledge is referred to a) Non-trivial extraction of implicit previously unknown and potentially useful information from data b) Set of columns in a database table that can be used to identify each record within this table uniquely c) collection of interesting and useful patterns in a database d) none of these	CO1	C
	<u>MCQ</u>	A star schema has what type of relationship between a dimension and fact table? a) Many-to-many b) One-to-one c) One-to-many d) All of the above	CO6	C
	<u>MCQ</u>	A Snowflake schema is which of the following types of table? a) Fact b) Dimension c) Helper d) All of the above	CO6	D
	<u>MCQ</u>	Bias is a) A class of learning algorithm that tries to find an optimum classification of a set of examples using the probabilistic theory	CO6	B

		b) Any mechanism employed by a learning system to constrain the search space of a hypothesis c) An approach to the design of learning algorithms that is inspired by the fact that when people encounter new situations, they often explain them by reference to familiar experiences, adapting the explanations to fit the new situation. d) None		
<u>Q.No:7</u>	<u>SAT</u>	Differentiate Between classification and clustering?	CO4	
	<u>SAT</u>	What is the difference between discrimination and classification?	CO4	
	<u>SAT</u>	Differentiate Between classification and regression?	CO4	
	<u>SAT</u>	Differentiate between metadata and data mart.	CO2	

SECTION-B

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

Time: 1 Hour and 30 Minutes

(3×12=36 Marks)

<u>Q No</u>	<u>Question</u>	<u>co</u>
<u>Q.No:8</u>	<p>Following data given for the attribute age: 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 30, 33, 33, 35, 35, 35, 35, 36, 40, 45, 46, 52, 70.</p> <p>(a) Use min-max normalization to transform the value 35 for age onto the range [0.0, 1.0].</p> <p>(b) Use z-score normalization to transform the value 35 for age, where the standard deviation of age is 12.94 years.</p> <p>(c) Use normalization by decimal scaling to transform the value 35 for age.</p> <p>Discuss on Normalization? Explain min-max normalization, z-score normalization and decimal scaling methods using appropriate example.</p> <p>What is Normalization? Use min-max normalization, z-score normalization and decimal scaling methods to normalize the following group of data: 200, 300, 400, 600, 1000</p>	CO2
<u>Q.No:9</u>	<p>Consider the following dataset represented by 5 training example. The target attribute is acceptable which can have values yes and no. Construct a decision tree from the given table. Show the value of information gain for each candidate attribute at each step in the construction of the tree.</p>	CO4

	<table><tr><th>House</th><th>Furniture</th><th>Nos. rooms</th><th>New kitchen</th><th>Acceptable</th></tr><tr><td>1</td><td>No</td><td>3</td><td>Yes</td><td>Yes</td></tr><tr><td>2</td><td>Yes</td><td>3</td><td>No</td><td>No</td></tr><tr><td>3</td><td>No</td><td>4</td><td>No</td><td>Yes</td></tr><tr><td>4</td><td>No</td><td>3</td><td>No</td><td>No</td></tr><tr><td>5</td><td>Yes</td><td>4</td><td>No</td><td>Yes</td></tr></table>	House	Furniture	Nos. rooms	New kitchen	Acceptable	1	No	3	Yes	Yes	2	Yes	3	No	No	3	No	4	No	Yes	4	No	3	No	No	5	Yes	4	No	Yes																
House	Furniture	Nos. rooms	New kitchen	Acceptable																																											
1	No	3	Yes	Yes																																											
2	Yes	3	No	No																																											
3	No	4	No	Yes																																											
4	No	3	No	No																																											
5	Yes	4	No	Yes																																											
<p>Consider the following transactional database T. Let min sup = 60% and min conf = 80%.</p> <table><tr><th>TID</th><th>Items bought</th></tr><tr><td>T100</td><td>{M, O, N, K, E, Y}</td></tr><tr><td>T200</td><td>{D, O, N, K, E, Y }</td></tr><tr><td>T300</td><td>{M, A, K, E}</td></tr><tr><td>T400</td><td>{M, U, C, K, Y}</td></tr><tr><td>T500</td><td>{C, O, O, K, I, E}</td></tr></table> <p>a) Find all frequent itemsets using Apriorialgorithms.</p> <p>b) Which of the itemsets from a) are closed? Which of the itemsets from a) are maximal?</p> <p>c) Determine strong association rules.</p>			TID	Items bought	T100	{M, O, N, K, E, Y}	T200	{D, O, N, K, E, Y }	T300	{M, A, K, E}	T400	{M, U, C, K, Y}	T500	{C, O, O, K, I, E}																																	
TID	Items bought																																														
T100	{M, O, N, K, E, Y}																																														
T200	{D, O, N, K, E, Y }																																														
T300	{M, A, K, E}																																														
T400	{M, U, C, K, Y}																																														
T500	{C, O, O, K, I, E}																																														
<p>A simple example from the stock market involving only discrete ranges has Profit as categorical attributes, with values (up, down) and the training data is,</p> <table><tr><th>Age</th><th>Competitio</th><th>Type</th><th>Profit</th></tr><tr><td>Old</td><td>Yes</td><td>Software</td><td>Down</td></tr><tr><td>Old</td><td>No</td><td>Software</td><td>Down</td></tr><tr><td>Old</td><td>No</td><td>Hardware</td><td>Down</td></tr><tr><td>Mid</td><td>Yes</td><td>Software</td><td>Down</td></tr><tr><td>Mid</td><td>Yes</td><td>Hardware</td><td>Down</td></tr><tr><td>Mid</td><td>No</td><td>Hardware</td><td>Up</td></tr><tr><td>Mid</td><td>No</td><td>Software</td><td>Up</td></tr><tr><td>New</td><td>Yes</td><td>Software</td><td>Up</td></tr><tr><td>New</td><td>No</td><td>Hardware</td><td>Up</td></tr><tr><td>New</td><td>No</td><td>Software</td><td>Up</td></tr></table> <p>Apply the decision tree algorithm and show the generated rules.</p>			Age	Competitio	Type	Profit	Old	Yes	Software	Down	Old	No	Software	Down	Old	No	Hardware	Down	Mid	Yes	Software	Down	Mid	Yes	Hardware	Down	Mid	No	Hardware	Up	Mid	No	Software	Up	New	Yes	Software	Up	New	No	Hardware	Up	New	No	Software	Up	
Age	Competitio	Type	Profit																																												
Old	Yes	Software	Down																																												
Old	No	Software	Down																																												
Old	No	Hardware	Down																																												
Mid	Yes	Software	Down																																												
Mid	Yes	Hardware	Down																																												
Mid	No	Hardware	Up																																												
Mid	No	Software	Up																																												
New	Yes	Software	Up																																												
New	No	Hardware	Up																																												
New	No	Software	Up																																												
Q.No:10	<p>What is data mining? Describe the steps involved in data mining when viewed as a process of knowledge discovery.</p> <p>What is data pre-processing. Explain various data pre-processing methods used in data mining.</p> <p>Describe data mining from “ Business Intelligence” perspective. What are the various application in data mining.</p>	CO1																																													
Q.No:11	<p>Explain Hierarchical method clustering of classification with example?Construct the single link agglomerative hierarchical clustering for the given distance matrix:</p>	CO5																																													

	1	2	3	4	5
1	0				
2	2	0			
3	6	3	0		
4	10	9	7	0	
5	9	8	5	4	0

Consider the following data set consisting of the scores of two variables on each of seven individuals.

Subject	A	B
1	1.0	1.0
2	1.5	2.0
3	3.0	4.0
4	5.0	7.0
5	3.5	5.0
6	4.5	5.0
7	3.5	4.5

Apply k-Means algorithm to this data set and grouped into two clusters. As a first step in finding a sensible initial partition, let the A & B values of the two individuals furthest apart (using the Euclidean distance measure), define the initial cluster means, giving;

	Individual	Mean Vector (centroid)
Group 1	1	(1.0, 1.0)
Group 2	4	(5.0, 7.0)

Consider five points {X1, X2, X3, X4, X5} with the following coordinates as a 2D sample for clustering : X1 =(0,0.25), X2=(0, 0), X3=(1.5, 0), X4=(5, 0), X5=(5, 2). Illustrate K-Means partitioning algorithm using the given dataset for two cluster.