



Sample Question Format

(For all courses having end semester Full Mark=50)

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

Subject Name & Code: Data Mining and Data Warehousing (IT-3031)

Applicable to Courses:

B.Tech

Full Marks=50

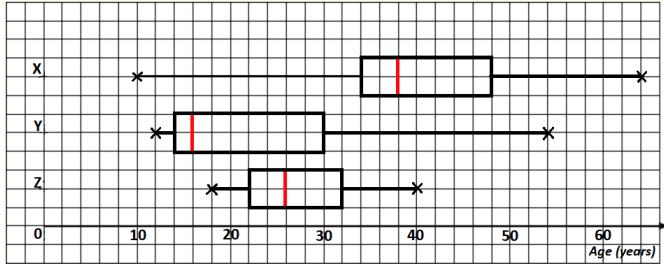
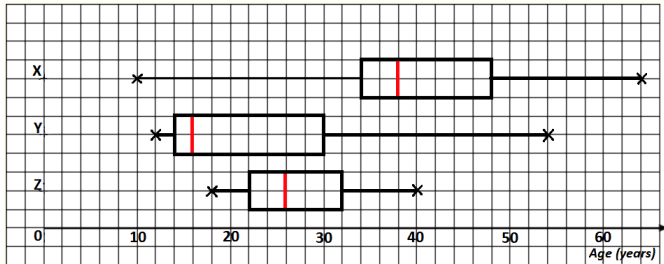
Time:2 Hours

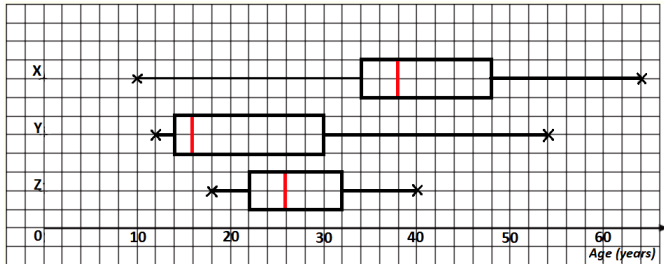
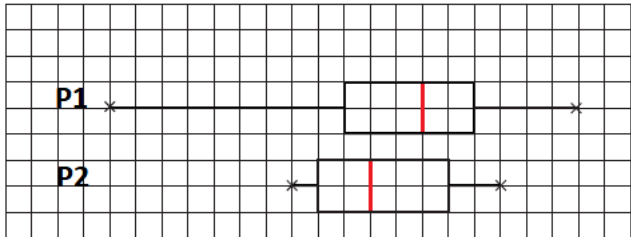
SECTION-A(Answer All Questions. Each question carries 2 Marks)

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>	Major data mining activities include the following general operations except, A. Exploratory data analysis B. Predictive modeling C. Discovering patterns and rules D. Data interpretation	CO1	D
	<u>MCQ</u>	Assumption is satisfied when the probability of missing values in one variable is unrelated to the value of the variable itself or to values of any other variable, A. Assumption of Missing at Random B. Assumption of Missing Completely at Random C. Assumption of Missing Not at Random D. Assumption of Missing Completely not at Random	CO1	B
	<u>MCQ</u>	What are the difficulties in implementing a data warehouse? (I). Construction, (II). Administration, (III). Quality control, (IV). Building blocks A. I, II, III B. I, II, IV	CO2	A

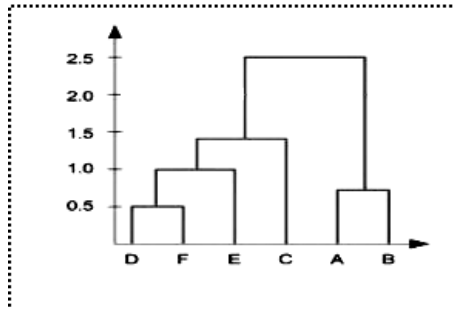
		C. II, III, IV D. I, III, IV		
	MCQ	Which one of the following is helps to store the data closer to users to enhance the performance. A. Data warehouse B. Data mart C. Metadata D. None	CO2	B
Q.No: 2	MCQ	<p>Three different games are playing by three differnt age group (X, Y, and Z) of the people. Players behaviours are visualizing in the following box ploat with whisker, answer the question;</p>  <p>Which game do you think <i>you</i> (according to your age) would not be allowed to play?</p> <p>A. Game X B. Game Y C. Game Z D. None</p>	CO1/ CO2	C
	MCQ	<p>Three different games are playing by three differnt age group (X, Y, and Z) of the people. Players behaviours are visualizing in the following box ploat with whisker, answer the question;</p>  <p>Which game would <i>you</i> probably enjoy most?</p> <p>A. Game X B. Game Y C. Game Z D. None</p>	CO1/ CO2	B
	MCQ	Three different games are playing by three differnt age group (X, Y, and Z) of the people. Players behaviours are visualizing in the following box ploat with whisker, answer the question;	CO1/ CO2	A

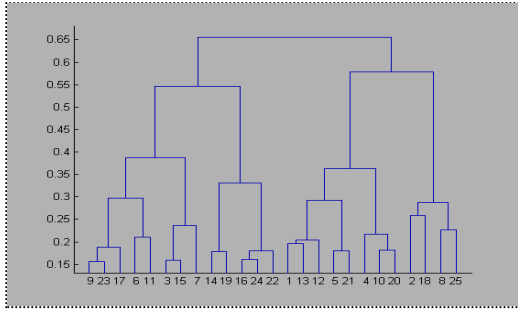
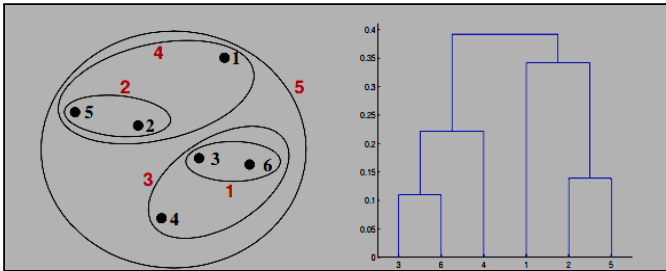
		 <p>Which game would <i>your parents</i> probably enjoy most?</p> <p>A. Game X B. Game Y C. Game Z D. None</p>		
	MCQ	<p>In an industry, the quality inspector will check the two different types of product that are marked in P1 & P2. The box and whisker plots below shows the results of quality tests on the quality of the respective products.</p>  <p>“The inspector might prefer to use the product (P2) because that P2 has the smaller _____, and the larger _____, which means that industry is less likely to produce a poor product.</p> <p>A. range, minimum B. range, maximum C. median, range D. lower quartile, median</p>	CO1/ CO2	A
Q.No: 3	MCQ	<p>In a shopping mall the products are numbered from 01 to 30 are Sampos and numbered from 31 to 50 are Conditioners. Which product numbers would you include in a systematic sample of size 10 ?</p> <p>A. 10, 20, 30, 40, 50 B. 01, 06, 11, 16, 21, 26, 31, 36, 41, 46 C. 01, 11, 21, 31, 41 D. 05, 10, 15, 20, 25, 30, 35, 40, 45, 50</p>	CO2/C O3	D
	MCQ	<p>Our KIIT University employs the following numbers of faculty members in 3 different positions as; Professor: 10, Associate Prof.: 20, Asst. Prof: 20. How many from each positions should be included in a quota sample of size: 25?</p> <p>A. 3, 11, 11 B. 5, 10, 10</p>	CO1/ CO2	B

		C. 10, 5, 10 D. 10, 10, 5																		
	MCQ	In a picnic trip consists of 40 members of whom 15 are gents. A quota of size 8 is to be selected for site visit. How many ladies and how many gents should be included in the sample? A. 6, 2 B. 5, 3 C. 2, 6 D. 3, 5	CO1/ CO2/ CO3	B																
	MCQ	In a DMDW class 75 students are present of of whom 15 are girls. A quota of size 15 is to be selected for site visit. How many boys and how many girls should be included in the sample? A. 5, 10 B. 10, 5 C. 3, 12 D. 12, 3	CO1	D																
Q.No: 4	MCQ	The ring sizes for the customers of a jewellery shop are shown in the table below. What will be the mean value? <table border="1"><thead><tr><th>Waist size</th><th>Frequency</th></tr></thead><tbody><tr><td>4</td><td>2</td></tr><tr><td>5</td><td>4</td></tr><tr><td>6</td><td>7</td></tr><tr><td>7</td><td>5</td></tr><tr><td>8</td><td>6</td></tr><tr><td>9</td><td>3</td></tr><tr><td>10</td><td>3</td></tr></tbody></table> A. 7 B. 4.28 C. 30 D. 1.63	Waist size	Frequency	4	2	5	4	6	7	7	5	8	6	9	3	10	3	CO1/ CO2	A
Waist size	Frequency																			
4	2																			
5	4																			
6	7																			
7	5																			
8	6																			
9	3																			
10	3																			
	MCQ	The number of minutes in a bus stop for a particular bus service was late have been shown in the table as; <table border="1"><thead><tr><th>Minutes Late</th><th>Frequency</th></tr></thead><tbody><tr><td>on time</td><td>19</td></tr><tr><td>1-5</td><td>12</td></tr><tr><td>6-10</td><td>9</td></tr><tr><td>11-20</td><td>4</td></tr><tr><td>21-40</td><td>4</td></tr><tr><td>41-60</td><td>2</td></tr><tr><td>over 60</td><td>0</td></tr></tbody></table> Estimate the probability of a bus being more than 20 minutes late. A. 8%	Minutes Late	Frequency	on time	19	1-5	12	6-10	9	11-20	4	21-40	4	41-60	2	over 60	0	CO1/ CO2	B
Minutes Late	Frequency																			
on time	19																			
1-5	12																			
6-10	9																			
11-20	4																			
21-40	4																			
41-60	2																			
over 60	0																			

		B. 12% C. 80% D. 88%																		
	MCQ	<p>The different items of a customer’s basket in a shopping mall are shown in the table below. What will be the mean value?</p> <table><tr><td>Item Codes</td><td>Frequency</td></tr><tr><td>3</td><td>2</td></tr><tr><td>4</td><td>4</td></tr><tr><td>5</td><td>7</td></tr><tr><td>6</td><td>5</td></tr><tr><td>7</td><td>6</td></tr><tr><td>8</td><td>3</td></tr><tr><td>9</td><td>3</td></tr></table> <p>A. 4 B. 6 C. 5 D. 1</p>	Item Codes	Frequency	3	2	4	4	5	7	6	5	7	6	8	3	9	3	CO1/ CO2	B
Item Codes	Frequency																			
3	2																			
4	4																			
5	7																			
6	5																			
7	6																			
8	3																			
9	3																			
	MCQ	<p>The number of minutes in a bus stop for a particular bus service was late have been shown in the table as;</p> <table><tr><td>Minutes Late</td><td>Frequency</td></tr><tr><td>on time</td><td>15</td></tr><tr><td>1-5</td><td>10</td></tr><tr><td>6-10</td><td>11</td></tr><tr><td>11-20</td><td>8</td></tr><tr><td>21-40</td><td>4</td></tr><tr><td>41-60</td><td>2</td></tr><tr><td>over 60</td><td>0</td></tr></table> <p>Estimate the probability of a bus being late of 10 minutes or less. A. 42% B. 72% C. 28% D. 88%</p>	Minutes Late	Frequency	on time	15	1-5	10	6-10	11	11-20	8	21-40	4	41-60	2	over 60	0	CO1/ CO2	A
Minutes Late	Frequency																			
on time	15																			
1-5	10																			
6-10	11																			
11-20	8																			
21-40	4																			
41-60	2																			
over 60	0																			
Q.No: 5	MCQ	<p>An itemset {Bread, Milk, Butter} whose support value is $10 \geq$ a minimum support threshold is considered as,</p> <p>A. Itemset B. Frequent Itemset C. Infrequent items D. Threshold values</p>	CO3	B																
	MCQ	<p>How do you calculate Confidence (Shoes→Socks), if support of Shoes is 5, support of Socks is 12, and support of together purchased is 60?</p> <p>A. 5 B. 12</p>	CO3	B																

		C. 1 D. 25														
	MCQ	In Association Rule Mining, which combination is correct? I.Support is never be equal to its confidence II.Support is always equal to its confidence III.Support is always greater than its confidence. IV. Support is always less than its confidence A. I, II, & III B. II, III, & IV C. ALL D. None	CO3	D												
	MCQ	Which of the following assumptions will satisfy, if the same minimum value of support is maintained at each level of an Association Rule Mining? I. An itemset is not frequent, then none of its supersets can be frequent. II. An itemset is frequent, then all its supersets are also frequent. III. An itemset is not frequent, then none of its subsets can be frequent. IV. An itemset is frequent, then all its subsets are also frequent. A. I & II B. II & III C. I & IV D. III & IV	CO3/ CO4	C												
Q.No: 6	MCQ	In a classification model if the model is predicted true for a class value whose actual value was false. Then this is a ____. A. False positive B. False negative C. True positive D. True negative	CO4	A												
	MCQ	Which one is incorrect option for support and confidence value for the following transaction data ? <table><tr><th>TID</th><th>ITEMS</th></tr><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></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/ CO4	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															
	MCQ	Perform KNN for "K=3" on the following dataset and	CO4	A												

		<p>generate the class level for the input (Acid durability =3 , strength=7, class=?).</p> <table><tr><th>Name</th><th>Acid durability</th><th>strength</th><th>class</th></tr><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></table> <p>A. Good B. Bad C. Invalid D. None</p>	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		
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																					
	MCQ	<p>In a brute-force approach for mining association rules, the total number of possible rules extracted from a data set that contains 4 items is,</p> <p>A. 40 B. 50 C. 51 D. 41</p>	CO5	B																				
Q.No: 7	MCQ	<p>What will be the no of clusters present in the below dendrogram?</p>  <p>A. 4 B. 5 C. 2 D. 3</p>	CO6	C																				
	MCQ	<p>At first iteration the 3 cluster observations using K-means algorithm are: C1: {(3,3), (5,5), (7,7)}; C2: {(0,4), (4,0), (8,8)}; C3: {(0,6), (6,0)}, What will be the cluster centroids if you want to proceed for second iteration?</p> <p>A. C1: (5,5), C2: (4,0), C3: (6,0) B. C1: (5,5), C2: (4,4), C3: (3,3) C. C1: (5,5), C2: (0,4), C3: (0,6) D. C1: (5,5), C2: (4,0), C3: (6,6)</p>	CO4	B																				
	MCQ	<p>What will be the possible no of clusters present in the below dendrogram?</p>	CO5/ CO6	B																				

		 <p>A. 2 B. 4 C. 7 D. 10</p>		
	MCQ	<p>From the following clustering representations and dendrogram, identify the type of link proximity function in hierarchical clustering.</p>  <p>A. MAX or Complete link B. MIN or Single link C. Average link D. None</p>	CO5/ CO6	A

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 (Each question should be from the same CO(s))</u>

Q.No:8	Find the group mean, median, & mode of the following data scored by students. <div><table><tr><th>Scores</th><th>Frequency</th></tr><tr><td>1-10</td><td>6</td></tr><tr><td>11-20</td><td>9</td></tr><tr><td>21 - 30</td><td>11</td></tr><tr><td>31 - 40</td><td>32</td></tr><tr><td>41 - 50</td><td>17</td></tr><tr><td>51 - 60</td><td>22</td></tr><tr><td>61 - 70</td><td>27</td></tr><tr><td>71 - 80</td><td>15</td></tr><tr><td>81-90</td><td>2</td></tr><tr><td>91 - 100</td><td>3</td></tr></table></div>	Scores	Frequency	1-10	6	11-20	9	21 - 30	11	31 - 40	32	41 - 50	17	51 - 60	22	61 - 70	27	71 - 80	15	81-90	2	91 - 100	3	CO1/CO2
	Scores	Frequency																						
	1-10	6																						
11-20	9																							
21 - 30	11																							
31 - 40	32																							
41 - 50	17																							
51 - 60	22																							
61 - 70	27																							
71 - 80	15																							
81-90	2																							
91 - 100	3																							
	As per the data list given below, prepare the partition into 4 bins using Equi-depth binning method and perform the smoothing by bin mean, bin median, and bin boundaries. 13,15,15,17,17,18,21,22,22,22,23,23,24,25,26,32,42,47,47,47,73,74,75,77																							
	The ages of the 112 people who admitted in a hospital are grouped as follows: <div><table><tr><th>Age</th><th>Number</th></tr><tr><td>0 - 9</td><td>20</td></tr><tr><td>10 - 19</td><td>21</td></tr><tr><td>20 - 29</td><td>23</td></tr><tr><td>30 - 39</td><td>16</td></tr><tr><td>40 - 49</td><td>11</td></tr><tr><td>50 - 59</td><td>10</td></tr><tr><td>60 - 69</td><td>7</td></tr><tr><td>70 - 79</td><td>3</td></tr><tr><td>80 - 89</td><td>1</td></tr></table></div> Calculate the mean, median & mode of these grouped data according to the above data.	Age	Number	0 - 9	20	10 - 19	21	20 - 29	23	30 - 39	16	40 - 49	11	50 - 59	10	60 - 69	7	70 - 79	3	80 - 89	1			
Age	Number																							
0 - 9	20																							
10 - 19	21																							
20 - 29	23																							
30 - 39	16																							
40 - 49	11																							
50 - 59	10																							
60 - 69	7																							
70 - 79	3																							
80 - 89	1																							
Q.No:9	1. The given contingency table summarizes supermarket transaction data, where <i>Bread</i> refers to the transactions containing <i>Bread</i> , and \overline{Bread} refers to the transactions that do not contain <i>Bread</i> . Similarly, <i>Butter</i> refers to the transactions containing <i>Butter</i> , and \overline{Butter} refers to the transactions that do not contain <i>Butter</i> . <div><table><tr><td></td><td><i>Bread</i></td><td>\overline{Bread}</td><td>\sum_{Row}</td></tr><tr><td><i>Butter</i></td><td>2000</td><td>500</td><td>2500</td></tr><tr><td>\overline{Butter}</td><td>1000</td><td>1500</td><td>2500</td></tr></table></div>		<i>Bread</i>	\overline{Bread}	\sum_{Row}	<i>Butter</i>	2000	500	2500	\overline{Butter}	1000	1500	2500	CO2/CO3										
	<i>Bread</i>	\overline{Bread}	\sum_{Row}																					
<i>Butter</i>	2000	500	2500																					
\overline{Butter}	1000	1500	2500																					

\sum_{Col}	3000	2000	5000																																								
<p>(a) Suppose that the association rule “Bread→Butter” is mined. Given a minimum support threshold of 25% and a minimum confidence threshold of 50%, is this association rule strong?</p> <p>(b) Based on the given data, is the purchase of Bread independent of the purchase of Butter? If not, what kind of <i>correlation</i> relationship exists between the two?</p> <p>Explain what is meant by association rule mining? For the table perform apriori algorithm,</p> <p>i. Determine the k-itemsets (frequent) obtained.</p> <p>ii. Justify the strong association rule that has been determined.</p> <table><tr><th>TID</th><th>Items</th></tr><tr><td>01</td><td>1, 3, 4, 6</td></tr><tr><td>02</td><td>2, 3, 5, 7</td></tr><tr><td>03</td><td>1, 2, 3, 5, 8</td></tr><tr><td>04</td><td>2, 5, 9, 10</td></tr><tr><td>05</td><td>1, 4</td></tr></table> <p>Assume min sup=30% and min conf=75%</p> <p>Given the following transactional database:</p> <table><tr><th colspan="4">Transactions</th></tr><tr><td>1</td><td>C</td><td>B</td><td>H</td></tr><tr><td>2</td><td>B</td><td>F</td><td>S</td></tr><tr><td>3</td><td>A</td><td>F</td><td>G</td></tr><tr><td>4</td><td>C</td><td>B</td><td>H</td></tr><tr><td>5</td><td>B</td><td>F</td><td>G</td></tr><tr><td>6</td><td>B</td><td>E</td><td>O</td></tr></table> <p>(a) We want to mine all the frequent itemsets in the data using the Apriori algorithm. Assume the minimum support level is 30%.</p> <p>(b) Find all the association rules that involve only B, C, H (in either left or right hand side of the rule). The minimum confi dence is 70%.</p>				TID	Items	01	1, 3, 4, 6	02	2, 3, 5, 7	03	1, 2, 3, 5, 8	04	2, 5, 9, 10	05	1, 4	Transactions				1	C	B	H	2	B	F	S	3	A	F	G	4	C	B	H	5	B	F	G	6	B	E	O
TID	Items																																										
01	1, 3, 4, 6																																										
02	2, 3, 5, 7																																										
03	1, 2, 3, 5, 8																																										
04	2, 5, 9, 10																																										
05	1, 4																																										
Transactions																																											
1	C	B	H																																								
2	B	F	S																																								
3	A	F	G																																								
4	C	B	H																																								
5	B	F	G																																								
6	B	E	O																																								
Q.No:10	<p>Consider the following data table where ”Sale” is a class attribute. The training data is shown below.</p> <table><tr><th>Price</th><th>Feedback</th><th>Warranty</th><th>Sale</th></tr><tr><td>LOW</td><td>NO</td><td>NO</td><td>YES</td></tr><tr><td>LOW</td><td>NO</td><td>YES</td><td>YES</td></tr><tr><td>HIGH</td><td>YES</td><td>NO</td><td>YES</td></tr><tr><td>HIGH</td><td>YES</td><td>YES</td><td>NO</td></tr><tr><td>LOW</td><td>YES</td><td>NO</td><td>NO</td></tr><tr><td>LOW</td><td>YES</td><td>YES</td><td>YES</td></tr><tr><td>HIGH</td><td>NO</td><td>NO</td><td>NO</td></tr></table> <p>Build a conditional probability table. Show how Naïve Bayesian method is used to classify the following test data</p> <table><tr><th>Price</th><th>Feedback</th><th>Warranty</th></tr><tr><td>HIGH</td><td>YES</td><td>NO</td></tr></table>			Price	Feedback	Warranty	Sale	LOW	NO	NO	YES	LOW	NO	YES	YES	HIGH	YES	NO	YES	HIGH	YES	YES	NO	LOW	YES	NO	NO	LOW	YES	YES	YES	HIGH	NO	NO	NO	Price	Feedback	Warranty	HIGH	YES	NO	CO1/ CO4/ CO5	
Price	Feedback	Warranty	Sale																																								
LOW	NO	NO	YES																																								
LOW	NO	YES	YES																																								
HIGH	YES	NO	YES																																								
HIGH	YES	YES	NO																																								
LOW	YES	NO	NO																																								
LOW	YES	YES	YES																																								
HIGH	NO	NO	NO																																								
Price	Feedback	Warranty																																									
HIGH	YES	NO																																									

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,

Age	Competition	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

Apply the decision tree algorithm and show the generated rules.

The following data table shows the details of a second hand car sale company. Apply the ID3 decision tree algorithm and show the suitable generated rules.

Color	Type	Doors	Tires	Class
Red	SUV	2	Whitewall	+
Blue	Minivan	4	Whitewall	-
Green	Car	4	Whitewall	-
Red	Minivan	4	Blackwall	-
Green	Car	2	Blackwall	+
Green	SUV	4	Blackwall	-
Blue	SUV	2	Blackwall	-
Blue	Car	2	Whitewall	+
Red	SUV	2	Blackwall	-
Blue	Car	4	Blackwall	-
Green	SUV	4	Whitewall	+
Red	Car	2	Blackwall	+
Green	SUV	2	Blackwall	-
Green	Minivan	4	Whitewall	-

Q.No:11

1. To cluster the following 8 data examples representing locations with (x, y) coordinates & distributed into 3 clusters: A1=(3,11), A2=(3,6), A3=(9,5), A4=(6,9), A5=(8,6), A6=(7,5), A7=(2,3), A8=(5,10). Suppose that the initial seeds (centers of each cluster) are A1, A4 and A7.

Run the k-means algorithm for 3 epoch and show the clusters.

2. Discuss how Spatial data mining is societally important.

1. The following eight points representing locations with (x, y) coordinates. Initial cluster centers are: A1(4, 12), A4(7, 10) and A7(3, 4). The distance function between two points a = (x₁, y₁) and b = (x₂, y₂) is defined as $D(a, b) = |x_2 - x_1| + |y_2 - y_1|$

Points	x	y
A1	4	12
A2	4	7
A3	10	6
A4	7	10
A5	9	7
A6	8	6
A7	3	4
A8	6	11

Use K-Means Algorithm to find the three cluster centers after the second iteration.

2. Discuss how text data mining is societally important.

CO5/CO6

1. What is hierarchical clustering? Apply Agglomerative Hierarchical Clustering and draw single link and average link dendogram for the following distance matrix.

	A	B	C	D	E
A	0	2	6	10	9
B	2	0	3	9	8
C	6	3	0	7	5
D	10	9	7	0	4
E	9	8	5	4	0

2. Discuss how multimedia data mining is societally important.