

**GUJARAT TECHNOLOGICAL UNIVERSITY****M.C.A -IV<sup>th</sup> SEMESTER–EXAMINATION – MAY- 2012****Subject code: 640005****Date: 19/05/2012****Subject Name: Data Warehousing & Data Mining (DWDM)****Time: 10:30 am – 01:00 pm****Total Marks: 70****Instructions:**

1. Attempt all questions.
2. Make suitable assumptions wherever necessary.
3. Figures to the right indicate full marks.

<b>Q.1</b>	<b>(a)</b>	Define Data Warehouse. Briefly describe the following terms: (i) Subject-oriented; (ii) Integrated; (iii) Time-variant; (iv) Non-volatile; (v) Fact; (vi) Dimension	<b>07</b>
	<b>(b)</b>	What is meant by Concept Hierarchy and what is its purpose? Is it applicable for Fact or for Dimension?	<b>03</b>
	<b>(c)</b>	Data Mining is classified into two categories: (i) Descriptive Data Mining, and (ii) Predictive Data Mining. Write briefly the basic purpose of these two categories of Data Mining. To which of the two categories the following data mining tasks fall into: <ul style="list-style-type: none"> <li>• Concept Description</li> <li>• Detection of an outlier</li> </ul>	<b>04</b>
<b>Q.2</b>	<b>(a)</b>	What is meant by Supervised Learning and Unsupervised Learning? A data mining task may use either supervised learning or unsupervised learning or either of the two or none of the two types. Indicate against each data mining task, the learning type, i.e. supervised or unsupervised or either or none: <ul style="list-style-type: none"> <li>• Classification</li> <li>• Clustering</li> <li>• Characterization</li> <li>• Discrimination</li> <li>• Association Rule Mining</li> </ul>	<b>07</b>
	<b>(b)</b>	Suppose that a data warehouse of a hospital consists of three dimensions, namely time, doctor, and patient with the concept hierarchy as follows: Time: day, month, quarter, year Doctor: doctor, specialization (e.g. ophthalmologist, pediatrician, etc.) Patient: patient, category (e.g. outdoor, indoor) There are two measures, namely count and charge, where charge is the fee that a doctor charges a patient for a visit. (i) Draw a star schema for the above data warehouse. (ii) Starting with the base cuboid [day, doctor, patient], what specific OLAP operations should be performed in order to list the total fee collected by each doctor in the year 2010?	<b>07</b>
		<b>OR</b>	
	<b>(b)</b>	Using the data warehouse given in the above example (Q. 2 (b) main part), let us assume that the hospital's interest is to analyze the volume of patients and the revenue generated (through fee) for a group of doctors under each specialization and for each doctor on monthly and yearly basis. Which of the pre-computed cubes (list out all the relevant data cubes) are required for this task? Justify your	<b>07</b>

		answer. In this case, whether metadata has any role to play? If yes, briefly describe the role of metadata. If only base cube is available, what OLAP operations will be required to do the task?	
<b>Q.3</b>	<b>(a)</b>	<p>The cube definition statement has the following syntax  define cube &lt;cube_name&gt; [&lt;dimension_list&gt;] : &lt;measure_list&gt;</p> <p>The dimension definition statement has the following syntax  define dimension &lt;dimension_name&gt; as (&lt;attribute_or_subdimension_list&gt;)</p> <p>Using the above syntax, define the fact constellation schema in DMQL for the two ‘facts’ and the corresponding ‘dimensions’ given below.  First ‘fact’: sales with four dimensions, namely time, item, branch, and location with the two measures as (i) value_sold (i.e. sales value), (ii) quantity_sold  Second ‘fact’: shipping with five dimensions, namely time, item, shipper, from_location, to_location with the two measures as (i) amount_shipped, and (ii) quantity_shipped.  Concept hierarchy of dimensions is given below:  time: day, month, quarter, year; item: item_name, brand, type, supplier_type  branch: branch_name, branch_type; location: city, state, country  shipper: shipper_name, shipper_location, shipper_type  Note-1: any location whether of shipper or from_location or to_location will have same hierarchy as that of location.  Note-2: if there is a difficulty in defining the fact constellation, then define the schema for the two facts separately in DMQL.</p>	<b>07</b>
	<b>(b)</b>	What is Data Mart? Differentiate between Data Mart and Data Warehouse. Describe 3-tier architecture of Data Warehouse.	<b>07</b>
		<b>OR</b>	
<b>Q.3</b>	<b>(a)</b>	Write down the salient points differentiating OLAP against OLTP. What are the basic characteristics of ROLAP, MOLAP, and HOLAP?	<b>07</b>
	<b>(b)</b>	Data pre-processing includes (i) Data Cleaning, (ii) Data Integration, (iii) Data Transformation, (iv) Data Reduction. Write briefly the basic tasks done under each type of pre-processing.	<b>07</b>
<b>Q.4</b>	<b>(a)</b>	<p>A database has the following transactions. Let min_sup = 60% and min_conf = 80%</p> <p>TID items_bought (in the form of brand-item_category)</p> <p>T01 {Sunset-Milk, Dairyland-Cheese, Best-Bread}</p> <p>T02 {Goldfarm-Apple, Dairyland-Milk, Best-Cheese, Wonder-Bread, Tasty-Pie}</p> <p>T03 {Westcoast-Apple, Dairyland-Milk, Wonder-Bread, Tasty-Pie}</p> <p>T04 {Sunset-Milk, Dairyland-Cheese, Wonder-Bread}</p> <p>List the frequent k-itemset for the largest k at the granularity of item_category (e.g. item<sub>i</sub> could be “Milk”) for the following rule template:  For all X ∈ transaction, buys(X, item<sub>1</sub>) ^ buys(X, item<sub>2</sub>) =&gt; buys(X, item<sub>3</sub>) [s, c]</p> <p>Also list all of the strong association rules (with their support s and confidence c) containing the frequent k-itemset for the largest k.</p>	<b>07</b>
	<b>(b)</b>	<p>Suppose that the data mining task is to cluster the following 8 points (with (x, y) representing location) into three clusters.</p> <p>A<sub>1</sub>(2, 10), A<sub>2</sub>(2, 5), A<sub>3</sub>(8, 4), B<sub>1</sub>(5, 8), B<sub>2</sub>(7, 5), B<sub>3</sub>(6, 4), C<sub>1</sub>(1, 2), C<sub>2</sub>(4, 9).</p> <p>The distance function is Manhattan distance. Suppose initially we assign A<sub>1</sub>, B<sub>1</sub>, and C<sub>1</sub> as the center of each cluster, respectively. Use the k-means algorithm to</p>	<b>07</b>

		add two points only, i.e. $A_2$ and $A_3$ , in appropriate clusters and compute the new centers of the clusters. When a new point is added in a cluster, the new center is computed as follows: (New Center) = Mean of (Old Center) and the (New Point Added). In case of a conflict, choose the cluster for which $ x_1 - C_x $ and $ y_1 - C_y $ are closer to each other, where $(C_x, C_y)$ are the coordinates of the new center.																																																													
		<b>OR</b>																																																													
<b>Q.4</b>	<b>(a)</b>	A database has four transactions. Let $\text{min\_sup} = 60\%$ and $\text{min\_conf} = 80\%$ . TID    date            items_bought T100 18-Jan-2011    {A, B, D, K} T200 18-Jan-2011    {A, B, C, D, E} T300 19-Jan-2011    {A, B, C, E} T400 22-Jan-2011    {A, B, C, D} Find all frequent itemsets using Apriori algorithm. List all of the strong association rules (with support $s$ and confidence $c$ ) matching the following meta rule, where $X$ is a variable representing customers, and $\text{item}_i$ denotes variables representing items (e.g. “A”, “B”, etc.): For all $X \in \text{transaction}$ , $\text{buys}(X, \text{item}_1) \wedge \text{buys}(X, \text{item}_2) \Rightarrow \text{buys}(X, \text{item}_3) [s, c]$	<b>07</b>																																																												
	<b>(b)</b>	Consider the 8 points given above in Q. 4 (b) main part. Use the distance function as Manhattan distance, and initially take $A_1$ , $B_1$ , and $C_1$ as the center of each cluster, respectively. If the similarity threshold is given as $\text{distance} \leq 4$ , which points will get marked as outliers?  If the initial center of each cluster is assigned as $A_1$ , $A_3$ , and $C_1$ then which points will get marked as outliers?	<b>07</b>																																																												
<b>Q.5</b>	<b>(a)</b>	The following table consists of training data from an employee database. Instead of repeating a tuple having the same values for “department”, “status”, “age”, and “salary”, only one instance of the tuple is included with the repeat count mentioned under the column (i.e. attribute) “count”. <table border="1" style="margin-left: 20px;"> <thead> <tr> <th>department</th> <th>status</th> <th>age</th> <th>salary</th> <th>count</th> </tr> </thead> <tbody> <tr><td>sales</td><td>senior</td><td>31...35</td><td>E</td><td>30</td></tr> <tr><td>sales</td><td>junior</td><td>26...30</td><td>A</td><td>40</td></tr> <tr><td>sales</td><td>junior</td><td>31...35</td><td>B</td><td>40</td></tr> <tr><td>systems</td><td>senior</td><td>31...35</td><td>H</td><td>05</td></tr> <tr><td>systems</td><td>junior</td><td>26...30</td><td>E</td><td>03</td></tr> <tr><td>systems</td><td>junior</td><td>21...25</td><td>E</td><td>20</td></tr> <tr><td>systems</td><td>senior</td><td>41...45</td><td>H</td><td>03</td></tr> <tr><td>marketing</td><td>senior</td><td>36...40</td><td>E</td><td>10</td></tr> <tr><td>marketing</td><td>junior</td><td>31...35</td><td>D</td><td>04</td></tr> <tr><td>marketing</td><td>senior</td><td>46...50</td><td>E</td><td>04</td></tr> <tr><td>marketing</td><td>junior</td><td>26...30</td><td>C</td><td>06</td></tr> </tbody> </table> (i) Draw a decision tree taking “department” as the root node, and taking “status” at the next level. (ii) Derive the decision rules with the associated probabilities. For example, if all the tuples at the leaf node belong to only one class, the probability will be 1 (i.e. 100%). However, if, say, 80% of the tuples belong to a particular class, the probability will be 0.8 (i.e. 80%).	department	status	age	salary	count	sales	senior	31...35	E	30	sales	junior	26...30	A	40	sales	junior	31...35	B	40	systems	senior	31...35	H	05	systems	junior	26...30	E	03	systems	junior	21...25	E	20	systems	senior	41...45	H	03	marketing	senior	36...40	E	10	marketing	junior	31...35	D	04	marketing	senior	46...50	E	04	marketing	junior	26...30	C	06	<b>07</b>
department	status	age	salary	count																																																											
sales	senior	31...35	E	30																																																											
sales	junior	26...30	A	40																																																											
sales	junior	31...35	B	40																																																											
systems	senior	31...35	H	05																																																											
systems	junior	26...30	E	03																																																											
systems	junior	21...25	E	20																																																											
systems	senior	41...45	H	03																																																											
marketing	senior	36...40	E	10																																																											
marketing	junior	31...35	D	04																																																											
marketing	senior	46...50	E	04																																																											
marketing	junior	26...30	C	06																																																											
	<b>(b)</b>	Bayesian theorem is stated as follows: $P(H   X) = P(X   H) P(H) / P(X)$ Where $X$ is a data sample whose class label is unknown. $H$ is the hypothesis such as that the data sample $X$ belongs to a specified class $C$ .	<b>03</b>																																																												

		<p>The above equation can be re-written as:  <math>P(X   H) = P(H   X) P(X) / P(H)</math>          Can we use any one of the above-stated two equations in the Bayesian theorem?          Justify your answer.</p>																																																																																					
	(c)	Describe Hold-out method and k-fold cross-validation method for assessing classifier accuracy.	04																																																																																				
		<b>OR</b>																																																																																					
<b>Q.5</b>	(a)	<p>Training data tuples from a customer database are given below:</p> <table border="1"> <thead> <tr> <th>RID</th> <th>age</th> <th>income</th> <th>student</th> <th>credit_rating</th> <th>Class:buys_laptop</th> </tr> </thead> <tbody> <tr><td>1</td><td>&lt;=30</td><td>high</td><td>no</td><td>excellent</td><td>yes</td></tr> <tr><td>2</td><td>&lt;=30</td><td>medium</td><td>no</td><td>fair</td><td>yes</td></tr> <tr><td>3</td><td>&lt;=30</td><td>low</td><td>yes</td><td>fair</td><td>yes</td></tr> <tr><td>4</td><td>&lt;=30</td><td>low</td><td>no</td><td>fair</td><td>no</td></tr> <tr><td>5</td><td>31...40</td><td>high</td><td>no</td><td>fair</td><td>no</td></tr> <tr><td>6</td><td>31...40</td><td>medium</td><td>no</td><td>excellent</td><td>yes</td></tr> <tr><td>7</td><td>31...40</td><td>low</td><td>yes</td><td>fair</td><td>yes</td></tr> <tr><td>8</td><td>31...40</td><td>low</td><td>yes</td><td>excellent</td><td>no</td></tr> <tr><td>9</td><td>&gt;40</td><td>high</td><td>no</td><td>excellent</td><td>yes</td></tr> <tr><td>10</td><td>&gt;40</td><td>medium</td><td>no</td><td>fair</td><td>no</td></tr> <tr><td>11</td><td>&gt;40</td><td>low</td><td>no</td><td>fair</td><td>no</td></tr> <tr><td>12</td><td>&gt;40</td><td>medium</td><td>no</td><td>excellent</td><td>no</td></tr> <tr><td>13</td><td>&gt;40</td><td>low</td><td>no</td><td>fair</td><td>no</td></tr> </tbody> </table> <p>(i) Draw a decision tree taking “credit_rating” as the root node, and taking “age” at the next level.</p> <p>(ii) Derive the decision rules with the associated probabilities. For example, if all the tuples at the leaf node belong to only one class, the probability will be 1 (i.e. 100%). However, if, say, 80% of the tuples belong to a particular class, the probability will be 0.8 (i.e. 80%).</p>	RID	age	income	student	credit_rating	Class:buys_laptop	1	<=30	high	no	excellent	yes	2	<=30	medium	no	fair	yes	3	<=30	low	yes	fair	yes	4	<=30	low	no	fair	no	5	31...40	high	no	fair	no	6	31...40	medium	no	excellent	yes	7	31...40	low	yes	fair	yes	8	31...40	low	yes	excellent	no	9	>40	high	no	excellent	yes	10	>40	medium	no	fair	no	11	>40	low	no	fair	no	12	>40	medium	no	excellent	no	13	>40	low	no	fair	no	07
RID	age	income	student	credit_rating	Class:buys_laptop																																																																																		
1	<=30	high	no	excellent	yes																																																																																		
2	<=30	medium	no	fair	yes																																																																																		
3	<=30	low	yes	fair	yes																																																																																		
4	<=30	low	no	fair	no																																																																																		
5	31...40	high	no	fair	no																																																																																		
6	31...40	medium	no	excellent	yes																																																																																		
7	31...40	low	yes	fair	yes																																																																																		
8	31...40	low	yes	excellent	no																																																																																		
9	>40	high	no	excellent	yes																																																																																		
10	>40	medium	no	fair	no																																																																																		
11	>40	low	no	fair	no																																																																																		
12	>40	medium	no	excellent	no																																																																																		
13	>40	low	no	fair	no																																																																																		
	(b)	Let the number of (independent) attributes in a database be 5 and let there be 4 class labels. If we have to use multi-layer, feed-forward artificial neural network (ANN) with back-propagation, how many nodes (neurons) will be in the input layer, in the output layer, and in one hidden layer. In case any assumption and / or any rule-of-thumb is used, it should be stated clearly.	03																																																																																				
	(c)	Describe Bagging (or Bootstrap aggregation) and Boosting for increasing classifier accuracy.	04																																																																																				

\*\*\*\*\*