GUJARAT TECHNOLOGICAL UNIVERSITY

M.C.A -IVth 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.

Q.1	(a)	Define Data Warehouse. Briefly describe the following terms: (i) Subject- oriented; (ii) Integrated; (iii) Time-variant; (iv) Non-volatile; (v) Fact; (vi) Dimension	07					
	(b) What is meant by Concept Hierarchy and what is its purpose? Is it appli Fact or for Dimension?							
	(c)	 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 taskd fall into: Concept Description Detection of an outlier 	04					
Q.2	(a)	 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: Classification Clustering Characterization Discrimination Association Rule Mining 	07					
	(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? 	07					
		OR						
	(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	07					

		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?			
Q.3	(a)	 define cube <cube_name> [<dimension_list>] : <measure_list></measure_list></dimension_list></cube_name> The dimension definition statement has the following syntax define dimension <dimension_name> as (<attribute_or_subdimension_list>)</attribute_or_subdimension_list></dimension_name> 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. 			
		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.			
	(b)	What is Data Mart? Differentiate between Data Mart and Data Warehouse. Describe 3-tier architecture of Data Warehouse.	07		
		OR			
Q.3	(a)	Write down the salient points differentiating OLAP against OLTP. What are the basic characteristics of ROLAP, MOLAP, and HOLAP?	07		
	(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.	07		
Q.4	(a)	A database has the following transactions. Let min_sup = 60% and min_conf = 80% TID items_bought (in the form of brand-item_category) T01 {Sunset-Milk, Dairyland-Cheese, Best-Bread} T02 {Goldfarm-Apple, Dairyland-Milk, Best-Cheese, Wonder-Bread, Tasty-Pie} T03 {Westcoast-Apple, Dairyland-Milk, Wonder-Bread, Tasty-Pie} T04 {Sunset-Milk, Dairyland-Cheese, Wonder-Bread} List the frequent k-itemset for the largest k at the granularity of item_category (e.g. item _i could be "Milk") for the following rule template: For all X ϵ transaction, buys(X, item ₁) ^ buys(X, item ₂) => buys(X, item ₃) [s, c] Also list all of the strong association rules (with their support s and confidence c) containing the frequent k-itemset for the largest k.	07		
	(b)	Suppose that the data mining task is to cluster the following 8 points (with (x, y) representing location) into three clusters. A ₁ (2, 10), A ₂ (2, 5), A ₃ (8, 4), B ₁ (5, 8), B ₂ (7, 5), B ₃ (6, 4), C ₁ (1, 2), C ₂ (4, 9). The distance function is Manhattan distance. Suppose initially we assign A ₁ , B ₁ , and C ₁ as the center of each cluster, respectively. Use the k-means algorithm to	07		

		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.					
Q.4	(a)	OR A database has four transactions. Let min_sup = 60% and min_conf = 80%. TID date items bought	07				
		The function of the strong for the strong for the strong association rules (with support s and confidence c) matching the following meta rule, where X is a variable representing customers, and item _i denotes variables representing items (e.g. "A", "B", etc.): For all X ϵ transaction, buys(X, item ₁) ^ buys(X, item ₂) => buys(X, item ₃) [s, c]					
	 (b) Consider the 8 points given above in Q. 4 (b) main part. Use the distance fund as Manhattan distance, and initially take A₁, B₁, and C₁ as the center of cluster, respectively. If the similarity threshold is given as distance <= 4, we 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?					
Q.5	(a))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". 					
		systems junior 2630 E 03					
		systems junior 2125 E 20					
		systemssenior4145H03marketingsenior3640E10					
		marketing junior 3135 D 04 marketing senior 4650 E 04					
		marketing junior 2630 C 06					
		 (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%). 					
	(b)	Bayesian theorem is stated as follows:	03				
		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.	0.5				

		The above equation can be re-written as: P(X H) = P(H X) P(X) / P(H) Can we use any one of the above-stated two equations in the Bayesian theorem? Justify your answer.						
	(c)	Describe Hold-out method and k-fold cross-validation method for assessing classifier accuracy.						
				2		OR		
Q.5	(a)	RID 1 2 3 4 5 6 7 8 9 10 11 12 13 (i) (ii)	age <=30 <=30 <=30 <=30 3140 3140 3140 3140 >40 >40 >40 >40 >40 >40 >40 >40 >10 crive t all the th be 1 (i.e)	income high medium low low high medium low low high medium low decision the the next le the decision uples at the 5. 100%). H	student no no yes no no no yes yes no no no no no no ree takin, evel. n rules w e leaf nod	excellent fair fair fair fair excellent fair excellent fair fair excellent fair fair g "credit_rating		
	(b)	class la (ANN) layer, in	tibels. If with ba n the out	we have to ack-propaga tput layer, a	o use mul ation, ho and in on	ti-layer, feed-fo w many nodes	atabase be 5 and let there be 4 orward artificial neural network (neurons) will be in the input In case any assumption and / or y.	03
	(c)		be Bagg er accura		Bootstrap	aggregation)	and Boosting for increasing	04
