20IT604/PE02

Hal	l Ti	cket	t Nu	mb	er:		

III/IV B.Tech (Regular) DEGREE EXAMINATION

July/August,2023

Sixth Semester

Time: Three Hours

Answer question 1 compulsory. Answer one question from each unit. Information Technology

Data Warehousing and Data Mining

Maximum: 70 Marks

(14X1 = 14Marks) (4X14=56 Marks)

			CO	BL	Μ
1	a)	What is transactional data?	CO1	L1	
		Transactional data is a type of data processing that consists of executing a number of			
		transactions occurring concurrently-online banking, shopping, order entry etc			
	b)	What is meant by concept hierarchy?	CO1	L1	
		A concept hierarchy defines a sequence of mappings from a set of low-level concepts to			
		higher-level, more general concepts.			
	c)	Define data characterization.	CO1	L1	
		Data characterization is a summarization of the general characteristics or features of a			
		target class of data. The data corresponding to the user-specified class are typically			
		collected by a query.			
	d)	What are outliers?	CO1	L1	
		Outlier is a data object that deviates significantly from the rest of the data objects and			
		behaves in a different manner.			
	e)	What do you mean by tree pruning?	CO2	L1	
		Tree pruning attempts to identify and remove such branches, with the goal of improving			
		classification accuracy on unseen data. Pruning reduces the complexity of the final			
		classifier, and hence improves predictive accuracy by the reduction of overfitting.			
	f)	Define binning.	CO2	L2	
		Binning, also known as discretization or bucketing, is a data preprocessing technique used			
		in data mining. It involves dividing a continuous variable into a set of smaller intervals or			
		bins and replacing the original values with the corresponding bin labels.			
	g)	List any two methods to handle missing data.	CO2	L1	
	-	i)Ignoring the tuple			
		ii) Use a global constant to fill in the missing value			
		iii) Use the attribute mean to fill in the missing value			
	h)	What is supervised learning.	CO2	L1	
		Supervised learning is a machine learning approach that's defined by its use of labeled			
		datasets. These datasets are designed to train or "supervise" algorithms into classifying data			
		or predicting outcomes accurately.			
	i)	Define Pattern? How are they represented?	CO3	L2	
		Pattern mining concentrates on identifying rules that describe specific patterns within the			
		data. These are data points that occur more often in the dataset. Market-basket analysis,			
		which identifies items that typically occur together in purchase transactions.			
	j)	What is sampling.	CO3	L1	
		Sampling is a process in statistical analysis where researchers take a predetermined number			
		of observations from a larger population.			
	k)	What is partitioning method in clustering?	CO3	L2	
		This clustering method classifies the information into multiple groups based on the			
		characteristics and similarity of the data.			
	1)	Define Dendrogram?	CO4	L1	
		A dendrogram is a diagram that shows the hierarchical relationship between objects. The			
		main use of a dendrogram is to work out the best way to allocate objects to clusters.			
	m)	Define divisive clustering.	CO4	L1	
		The divisive clustering algorithm is a top-down clustering approach, initially, all the points			
		in the dataset belong to one cluster and The cluster splitting process repeats until.			
		avantually and now aluster contains only a single chiest			
		eventuary, each new cluster contains only a single object.	1		

	n)	Write the equations for mean and average distance.	CO4	L2	
		Mean distance :			
		Average distance :			
		Tu:4 T			
2	0)	<u>Unit-1</u> Explain the design and construction of Data warehouse	CO1	12	7M
2	<i>a)</i>	To design an effective data warehouse we need to understand and analyze business needs	COI		/ 101
		and construct a <i>business analysis framework</i> . The construction of a large and complex			
		information system can be viewed as the construction of a large and complex building for			
		which the owner, architect, and builder have different views. These views are combined to			
		form a complex framework that represents the top-down, business-driven, or owner's			
		perspective, as well as the bottom-up, builder-driven, or implementor's view of the			
		information system.			
		Four different views regarding the design of a data warehouse must be considered: the top-			
		down view, the data source view, the data warehouse view, and the business query view.			
		• The top-down view allows the selection of the relevant information necessary for the			
		data warehouse. This information matches the current and future business needs.			
		• The data source view exposes the information being captured, stored, and managed by			
		operational systems. This information may be documented at various levels of detail and			
		accuracy, from individual data source tables to integrated data source tables. Data			
		sources are often modeled by traditional data modelling techniques, such as the entity-			
		relationship model or CASE (computer-aided software engineering) tools.			
		• The data warehouse view includes fact tables and dimension tables. It represents the			
		information that is stored inside the data warehouse, including precalculated totals and counts, as well as information regarding the source, data, and time of origin, added to			
		counts, as well as information regarding the source, date, and time of origin, added to provide historical context			
		• Finally, the business guery view is the perspective of date in the date werehouse from			
		• Finally, the business query view is the perspective of data in the data warehouse from the viewpoint of the and user			
-	b)	On what kinds of data can the data mining can be performed?	CO1	13	7M
	0)	Explanation of any 3 kind of databases can be considered	001		/ 11/1
		• A number of different data repositories on which mining can be performed. In principle.			
		data mining should be applicable to any kind of data repository, as well as to transient			
		data, such as data streams. Thus the scope of our examination of data repositories will			
		include relational databases, data warehouses, transactional databases, advanced			
		database systems, flat files, data streams, and the World Wide Web. Advanced database			
		systems include object-relational databases and specific application-oriented databases,			
		such as spatial databases, time-series databases, text databases, and multimedia			
		databases. The challenges and techniques of mining may differ for each of the repository			
		systems.			
		(OD)			
2		(UK) Emploin Ston Schome in multidimensional database with a mut diamensional	CO1	1.2	714
3	a)	Explain Star Schema in multidimensional database with a heat diagram.	COI	L2	/ 1/1
		Diagram -2M			
		Diagram -2141			
		Star Schema: A star schema is a type of data modeling technique used in data			
		warehousing to represent data in a structured and intuitive way. In a star schema, data is			
		organized into a central fact table that contains the measures of interest, surrounded by			
		dimension tables that describe the attributes of the measures.			
		The fact table is a star scheme contains the measures or matrice that are of interest to the			
		user or organization. For example, in a sales data warehouse, the fact table might contain			
		sales revenue units sold and profit margins. Each record in the fact table represents a			
		specific event or transaction, such as a sale or order.			
		The dimension tables in a star schema contain the descriptive attributes of the measures			
		in the fact table. These attributes are used to slice and dice the data in the fact table,			
		anowing users to analyze the data from different perspectives. For example, in a sales			
		location			

			dimedi dimedi time_key day_of_the month quarter year branch_un branch_n branch_ty	ch n table	sales fact table time_key item_key branch_key location_key dollar_sold	din fitte fi	item mension table sm_key in_name and pe pplier_type location mension table tion_key et vince_or_state ntry			
	b)	How to identify the interestingness of a pattern.							L3	7M
		A pattern is interesting if it is (1) <i>easily understood</i> by humans, (2) <i>valid</i> on new or test data with some degree of <i>certainty</i> , (3) potentially <i>useful</i> , and (4) <i>novel</i> . A pattern is also interesting if it validates a hypothesis that the user <i>sought to confirm</i> . An interesting pattern represents knowledge . Several objective measures of pattern interestingness exist. These are based on the structure of discovered patterns and the statistics underlying them. An objective measure for association rules of the form $X \Rightarrow Y$ is rule support , representing the percentage of transactions from a transaction database that the given rule satisfies. This is taken to be the probability $P(X \cup Y)$, where $X \cup Y$ indicates that a transaction contains both X and Y, that is, the union of itemsets X and Y. Another objective measure for association. This is taken to be the conditional probability $P(Y X)$, that is, the probability that a transaction containing X also contains Y. More formally, support and confidence are defined as $\begin{array}{c} Support(X \Rightarrow Y) = P(X \cup Y) \\ Confidence(X \Rightarrow Y) = P(Y X) \end{array}$ In general, each interestingness measure is associated with a threshold, which may be controlled by the user. For example, rules that do not satisfy a confidence threshold of, say,								
					τ	J nit-II				
4	a)	Consider the dat	aset below	and draw a	decision tr	ree using in	formation gain .	CO2	L3	7M
		Index	А	В	C	D	E (Class)			
		1	4.8	3.4	1.9	0.2	positive			
		2	5	3	1.6	1.2	positive			
		3	5	3.4	1.6	0.2	positive			
		4	5.2	3.5	1.5	0.2	positive			
		5	5.2	3.4	1.4	0.2	positive			
		6	4.7	3.2	1.0	0.2	positive			
		7	4.0	3.1	1.0	0.2	positivo			
		9	7	3.4	1.3 4 7	0.4	negative			
		10	64	3.2	4.7	1.4	negative			
		11	69	3.2	49	1.5	negative			
		12	5.5	2.3	4	1.3	negative			
		12	6.5	2.8	4.6	1.5	negative			
		14	5.7	2.8	4.5	1.3	negative			
		15	6.3	3.3	4.7	1.6	negative			
		16	4.9	2.4	3.3	1	negative			
		Solution can be	considered	using infor	mation gain	n or Gini in	dex			

b)	Use these methods to normalize the following group of data: 200, 300, 400, 600,1000	CO2	L3	7M
	(i) min-max normalization by setting min = 0 and max =1			
	(ii) Z-score normalization			
	(a) Min-Max Normalization Data given :- 200, 300, 400, 600, 1000			
	Min = 200			
	Max = 1000			
	V = the respective value of the attribute			
	V1 = 200			
	V2 = 300			
	V3 = 400			
	V4 = 600			
	V5 = 1000			
	New max $= 1$			
	New $\min = 0$			
	V' = {V-minA/maxA-minA(new maxA-new minA)} + new min A			
	Therefore,			
	for 200,			
	$\min \max = (200-200(1-0)/1000-200) + 0$			
	=0			
	for 300			
	$\min \max = (300-200(1-0)/1000-200) + 0$			
	= 100/800			
	= 0.125			
	for 400			
	$\min \max = (400-200(1-0)/1000-200) + 0$			
	= 200/800			
	= 0.25			
	for 600			
	$\min \max = 600-200/100-200 \times (1-0)+0$			
	= 400/800			
	= 0.5			
	for 1000			
	min max = 1000-200/1000-200×(1-0)+0			

		= 1			
		(b) Z-score normalization			
		200,300,400,600,1000			
		Standard derivation = $\sqrt{\sum}(\text{every individual data - mean of data})^2/n$			
		Now,			
		mean value of following data = $(200+300+400+600+1000)/5$			
		= 2500/5			
		= 500			
		Therefore,			
		Standard derivation = $\sqrt{\{(200-500)^2 + (300-500)^2 + (400-500)^2 + (600-500)^2 + (1000-500)^2\}/5}$			
		$= \sqrt{\{(-300)^2 + (-200)^2 + (-100)^2 + (100)^2 + (500)^2\}/5}$			
		$= \sqrt{(90000 + 40000 + 10000 + 10000 + 250000)/5}$			
		$=\sqrt{400000/5}$			
		$=\sqrt{80000}$			
		= 282.8			
		Z score = 200-500/282.8			
		= -1.06			
		Z score = 300-500/282.8			
		= -0.7			
		Z score = 400-500/282.8			
		= -0.35			
		Z score = 600-500/282.8			
		= 0.35			
		Z score = 1000-500/282.8			
		= 1.78			
5		(OR)	CO1	12	714
5	a)	How tree pruning in decision tree induction is useful? Explain various methods for pruning decision trees. When decision trees are built, many of the branches may reflect noise or outliers in the training data.Tree pruning methods address this problem of overfittingthe data. Tree pruning attempts to identify and remove such branches, with the goal of improving classification accuracy on unseen data. Decision trees can suffer from repetition and replication, making them overwhelming to interpret. Repetition occurs when an attribute is repeatedly tested along a given branch of the tree.	02	L3	/M
		In replication, duplicate subtrees exist within the tree. These situations can impede the accuracy and comprehensibility of a decision tree.			
		Pruned trees			
		• These tend to be smaller and less complex and, thus, easier to comprehend.			

	1				
		• They are usually faster and better at correctly classifying independent test data than unpruned trees.			
		• Pruned trees tend to be more compact than their unpruned counterparts			
		There are two common approaches to tree pruning:			
		1. prepruning :			
		•In the pre-pruning approach, a tree is "pruned" by halting its construction early (e.g. by deciding not to further split or partition the subset of training tuples at a given node).			
		•When constructing a tree, measures such as statistical significance, information gain, Gini index, and so on can be used to assess the goodness of a split.			
		• If partitioning the tuples at a node would result in a split that falls below a pre specified threshold, then further partitioning of the given subset is halted.			
		•There are difficulties, however, in choosing an appropriate threshold.			
		•High thresholds could result in oversimplified trees, whereas low thresholds could result in very little simplification.			
		2. post pruning.			
		•The second and more common approach is post pruning, which removes subtrees from a "fully grown" tree.			
		•A subtree at a given node is pruned by removing its branches and replacing it with a leaf.			
		•The leaf is labeled with the most frequent class among the subtree being replaced.			
		•The cost complexity pruning algorithm used in CART is an example of the post pruning approach.			
		•The basic idea is that the simplest solution is preferred.			
		•Unlike cost complexity, pruning does not require an independent set of tuples.			
		•Post pruning leads to a more reliable tree.			
	b)	Suppose a group of 12 sales price records has been sorted as follows: 5, 10, 11, 13, 15, 35, 50, 55, 72, 92, 204, 215. Partition them into three bins by each of the following methods: (i) equal-frequency (equal-depth) partitioning (ii) equal-width partitioning	CO2	L3	7M
		 (i) Equal-frequency partitioning Partition the data into equal depthbins of depth 4 Bin 1: 5, 10, 11 13 Bin 2: 15, 35, 50, 55 Bin 3: 72, 92, 204, 215 			
		(ii) Equal -width partitioning			
		5)/3=70. We get			
		Bin 1: 5, 10, 11 13, 15, 35, 50, 55, 72 Bin 2: 92			
		Bin 3: 204, 215			
_	`		002	1.2	71 5
6	a)	"Strong Kules Are Not Necessarily Interesting". Justify.	CO3	L3	/M
		Whether or not a rule is interesting can be assessed either subjectively or objectively.			
		Ultimately, only the user can judge if a given rule is interesting, and this judgment, being subjective may differ from one user to another. However, objective interesting ended			
		measures based on the statistics — behind the data can be used as one step toward the			
		goal of weeding out uninteresting rules from presentation to the user.			

the fo	the support-confidence framework for association rules. This leads to correlation rules of the form								
	$A \rightarrow B$ [support, confidence. correlation].								
That i	That is, a correlation rule is measured not only by its support and confidence but also by the								
correl	ation b	etween itemsets A a	nd B. There are	e many differen	nt correlat	ion mea	sures from		
which	which to choose. In this section, we study various correlation measures to determine which								
would be good for mining large data sets.									
Consider the following transactional data for a commercial shop.								CO3	L3
,		TID		List of Items	with Ids				
		T1		i2, i4					
		T2		i1, i2, i5					
		T3		i2, i3					
		T4		11, 13					
		15 T6		11, 12, 14					
		T0 T7		i1_i3					
		Т, Т8		i1, i2. i3					
		T9		i1, i2, i3, i5					
Gener	Generate all the frequent itemsets using apriori algorithm. Consider the minimum support								
count									
т	Transactions List				_				
Т		List of Itoms IDs		1-item Set	s F	requen	icy		
T.	100	11 12 15		11		7			
	100	11, 12, 13	_	12					
T	200	12,14		13		5	1		
T. T.	200 300	12, 14 12, 13	_	13		5 4			
T T T	200 300 400	12, 14 12, 13 11, 12, 14	_	13 14 15		5 4 2			
Т: Т: Т4 Т!	200 300 400 500	 I2, I4 I2, I3 I1, I2, I4 I1, I3 		13 14 15		5 4 2			
T: T: T: T: T: T:	200 300 400 500 600	 I2, I4 I2, I3 I1, I2, I4 I1, I3 I2, I3 	F	13 14 15 Frequent 1-ite	em Sets	5 4 2 Freq	uency		
T: T: T- T: T: T: T:	200 300 400 500 600 700	 I2, I4 I2, I3 I1, I2, I4 I1, I3 I2, I3 I1, I3 	F	13 14 15 Frequent 1-ito 11	em Sets	5 4 2 Freq	uency 6		
ד: ד: די די די די	200 300 400 500 600 700 800	 I2, I4 I2, I3 I1, I2, I4 I1, I3 I2, I3 I1, I3 I1, I3 I1, I3 I1, I3, I5 	F	13 14 15 Frequent 1-ite 11 12	em Sets	5 4 2 Freq	uency 6 7		
די די די די די די די די די	200 300 400 500 600 700 800 900	 12, 14 12, 13 11, 12, 14 11, 13 12, 13 11, 13 11, 12, 13, 15 11, 12, 13 	F	13 14 15 Frequent 1-ite 11 12 13	em Sets	5 4 2 Freq	uency 6 7 5		
T; T; T; T; T; T; T; T; T;	200 300 400 500 600 700 800 900	 12, 14 12, 13 11, 12, 14 11, 13 12, 13 11, 13 11, 12, 13, 15 11, 12, 13 	F	13 14 15 Frequent 1-ite 11 12 13 14	em Sets	5 4 2 Freq	uency 6 7 5 4		
ד: ד: די די די די די	200 300 400 500 600 700 800 900	 I2, I4 I2, I3 I1, I2, I4 I1, I3 I2, I3 I1, I3 I1, I2, I3, I5 I1, I2, I3 	F	13 14 15 Frequent 1-ite 11 12 13 14	em Sets	5 4 2 Freq	uency 6 7 5 4		
ד: די די די די די די די	200 300 400 500 600 700 800 900	 12, 14 12, 13 11, 12, 14 11, 13 12, 13 11, 13 11, 12, 13, 15 11, 12, 13 	F	13 14 15 Frequent 1-ite 11 12 13 14	em Sets	5 4 2 Freq	uency 6 7 5 4		
Т; Т; Т; Т; Т; Т; Т; Т;	200 300 400 500 600 700 800 900	 12, 14 12, 13 11, 12, 14 11, 13 12, 13 11, 13 11, 12, 13, 15 11, 12, 13 		13 14 15 Frequent 1-ite 11 12 13 14	em Sets	5 4 2 Freq	uency 6 7 5 4		
נד די די די די די די די	200 300 400 500 600 700 800 900	 12, 14 12, 13 11, 12, 14 11, 13 12, 13 11, 13 11, 12, 13, 15 11, 12, 13 	F	13 14 15 Frequent 1-ite 11 12 13 14	em Sets	5 4 2 Freq	uency 6 7 5 4		
ד: ד: ד: ד: ד: ד: ד: ד:	200 300 400 500 600 700 800 900	 12, 14 12, 13 11, 12, 14 11, 13 12, 13 11, 13 11, 12, 13, 15 11, 12, 13 	2-item Sets	13 14 15 Frequent 1-ite 11 12 13 14	em Sets	5 4 2 Freq	uency 6 7 5 4		
T T T T T T T T T T T T T T T	200 300 400 500 600 700 800 900	 12, 14 12, 13 11, 12, 14 11, 13 12, 13 11, 13, 15 11, 12, 13, 15 11, 12, 13 	2-item Sets	13 14 15 Frequent 1-ite 11 12 13 14	em Sets	5 4 2 Freq	uency 6 7 5 4		
T; T; T; T; T; T; T; T; T; T; T; T; T; T	200 300 400 500 600 700 800 900	12, 14 12, 13 11, 12, 14 11, 13 12, 13 11, 12, 13, 15 11, 12, 13, 15 11, 12, 13	2-item Sets 11, 12	13 14 15 irequent 1-ite 11 12 13 14	em Sets	5 4 2 Freq	uency 6 7 5 4		
T; T; T; T; T; T; T; T; T; T; T; T; T; T	200 300 400 500 600 700 800 900 fransac	12, 14 12, 13 11, 12, 14 11, 13 12, 13 11, 13 11, 12, 13, 15 11, 12, 13	2-item Sets 11, 12 11, 13 11, 14	13 14 15 irequent 1-ite 11 12 13 14	Freque 2-item	5 4 2 Freq	uency 6 7 5 4		
T; T; T; T; T; T; T; T; T; T; T; T; T; T	200 300 400 500 600 700 800 900 900 900	12, 14 12, 13 11, 12, 14 11, 13 12, 13 11, 12, 13, 15 11, 12, 13, 15 11, 12, 13 11, 12, 13 11, 12, 13 11, 12, 13 11, 12, 13	2-item Sets 11, 12 11, 13 11, 14 11, 15	13 14 15 Frequent 1-ite 11 12 13 14	Freque 2-item	5 4 2 Freq 5 5 5 5 5 5 5 5 2	uency 6 7 5 4 4		
T; T; T; T; T; T; T; T; T; T; T; T; T; T	200 300 400 500 600 700 800 900 900 0 0 0 0 0 0 0	12, 14 12, 13 11, 12, 14 11, 13 12, 13 11, 12, 13, 15 11, 12, 13, 15 11, 12, 13 11, 12, 13 11, 12, 13 11, 12, 13 11, 12, 13 11, 12, 13 12, 13 11, 12, 13	2-item Sets 11, 12 11, 13 11, 14 11, 15 12, 13	13 14 15 irequent 1-ite 11 12 13 14	Freque 2-item	5 4 2 Freq 5 5 5 6 1 5 6 1 5 6 1 5 6 1 5 1 5 1 5 1	uency 6 7 5 4 4		
T; T; T; T; T; T; T; T; T; T; T; T; T; T	200 300 400 500 600 700 800 900 900 900 900 900 900 900 900 9	12, 14 12, 13 11, 12, 14 11, 13 12, 13 11, 13 11, 12, 13, 15 11, 12, 13, 15 11, 12, 13, 15 11, 12, 13 11, 12, 13 11, 12, 13 11, 12, 13 12, 13 11, 12, 15 12, 14 12, 13 11, 12, 14 11, 13	2-item Sets 11, 12 11, 13 11, 14 11, 15 12, 13 12, 14	13 14 15 irequent 1-ite 11 12 13 14 12 13 14	Freque 2-item 11, 11 11, 11	5 4 2 Freq 5 5 5	uency 6 7 5 4 4 requency 4 4 2		
T; T; T; T; T; T; T; T; T; T; T; T; T; T	200 300 400 500 600 700 800 900 900 0 0 0 0 0 0 0 0 0 0 0 0	 12, 14 12, 13 11, 12, 14 11, 13 12, 13 11, 12, 13, 15 11, 12, 13, 15 11, 12, 13 	2-item Sets 11, 12 11, 13 11, 14 11, 15 12, 13 12, 13 12, 14 12, 15	13 14 15 irequent 1-ite 11 12 13 14 12 13 14	Freque 2-item 11, 12 11, 11 12, 12	5 4 2 Freq 5 2 3 3 3 3 4 5 3 3 4	uency 6 7 5 4 4 requency 4 4 2 3		
T; T; T; T; T; T; T; T; T; T; T; T; T; T	200 300 400 500 600 700 800 900 900 0 0 0 0 0 0 0 0 0 0 0 0	 i2, i4 i2, i3 i1, i2, i4 i1, i3 i2, i3 i1, i2, i3, i5 i1, i2, i3, i5 i1, i2, i3 i1, i2, i3 i1, i2, i3 i1, i2, i3 i1, i2, i4 i2, i3 i1, i3 i3, i3 	2-item Sets 11, 12 11, 13 11, 14 11, 15 12, 13 12, 14 12, 15 13, 14	13 14 15 Frequent 1-its 11 12 13 14 Frequency 4 4 1 2 3 2 3 2 0	Freque 2-item 11, 11 11, 11 12, 14 12, 14	5 4 2 Freq 5 3 3 4 4	uency 6 7 5 4 4 4 4 2 3 2 2		
T; T; T; T; T; T; T; T; T; T; T; T; T; T	200 300 400 500 600 700 800 900 900 900 0 0 0 0 0 0 0 0 0 0 0 0	 12, 14 12, 13 11, 12, 14 11, 13 12, 13 11, 13, 15 11, 12, 13, 15 11, 12, 13 tions List t.ist of Items IDs 11, 12, 15 12, 14 12, 13 11, 12, 14 11, 13 12, 13 11, 12, 13, 15	2-item Sets 11, 12 11, 13 11, 14 11, 15 12, 13 12, 13 12, 14 12, 15 13, 14	13 14 15 requent 1-its 11 12 13 14 12 13 14 12 13 14 12 3 2 3 2 0 1	Freque 2-item 11, 11 11, 11 12, 11 12, 11 12, 11	5 4 2 Freq 5 2 3 3 4 5 5 3 4 5 5 3 4 5 5 3 4 5 5 3 4 5 5 3 4 5 5 3 5 3	uency 6 7 5 4 4 7 5 4 4 7 5 4 4 7 5 4 4 7 7 5 4 4 7 7 5 4 4 7 7 5 4 7 7 7 7		

7 a) Discuss FP Growth algorithm with an example. CO3 L2 7M 7 a) Discuss FP Growth algorithm is an alternative way to find frequent item sets without using candidate generations, thus improving performance. For so much, it uses a divide-and-conquery strategy. The core of this method is the uses get of a special data structure named frequent-pattern tree (FP-tree), which retains the item set association information. CO3 L2 7M This algorithm works as follows: 0 First, it compresses the input database creating an FP-tree instance to represent frequent items. 0 Alter this first step, it divides the compressed database into a set of conditional databases, each associated with one frequent pattern. 0 Finally, each such database is mined separately. Using this strategy, the FP-Growth reduces the search costs by recursively looking for short patterns and then concatenating them into the long frequent patterns Image: the frequency of each individual frequency frequenc													
The FP-Growth Algorithm is an alternative way to find frequent item sets without using condidate generations, thus improving performance. For so much, it uses a divide-and-crequent strategy. The core of this method is the usage of a special data structure named frequent-pattern tree (FP-tree), which retains the item set association information. This algorithm works as follows: • First, it compresses the input database creating an FP-tree instance to represent frequent items. • Alter this first step, it divides the compressed database into a set of conditional databases, each associated with one frequent pattern. • Finally, cach such database is mined separately. Using this strategy, the FP-Growth reduces the search costs by recursively looking for short patterns and then concatenating them into the long frequent patterns Transaction ID Items T1 {E,K,M,M,O,Y} T2 {D,E,K,N,O,Y} T3 {A,E,K,M} T4 {C,K,M,U,Y} T5 {C,E,I,K,O,O} Item is computed:- Items Item is computed:- Items Item is computed:- Items Item is computed:- Items Item is computed:- Item is computed:- Item is computed:- Item is computed:- Item is computed:- Item is computed:- Ite	7	a)	Discuss FP Growth algorithm	n with an exa	imple.			CO3	L2	7M			
This algorithm works as follows:•First, it compresses the input database creating an FP-tree instance to represent frequent items.•After this first step, it divides the compressed database into a set of conditional databases, each associated with one frequent pattern.•Finally, each such database is mined separately.Using this strategy, the FP-Growth reduces the search costs by recursively looking for short patterns and then concatenating them into the long frequent patternsTransaction IDItemsII{E,K,M,N,O,Y}IG,E,K,MQT3{A,E,K,M}item is computed:-ItemsI[E,K,M,N,O,Y]I1IEI[E,K,M,O,Y]I1I[E,K,M,O,Y]I1I[E,K,M,N,O,Y]I1I[E,K,M,N,O,Y]I1I[E,K,M,N,O,Y]I1I[E,K,M,N,O,Y]I1I[E,K,M,N,O,Y]I1I[E,K,M,0,Y]I1I[E,K,M,0,Y]I1I[E,K,M,0,Y]I1I[E,K,M,0,Y]I1I[E,K,M,0,Y]I1I[E,K,M,0,Y]I1I[E,K,M,0,Y]I1I[E,K,M,0,Y]I1I[E,K,M,0,Y] <td></td> <td></td> <td>The FP-Growth Algorithm is candidate generations, thus in conquer strategy. The core of frequent-pattern tree (FP-tree)</td> <td>an alternative mproving per this method which retains</td> <td>e way to find fre formance. For s is the usage of a s the item set asso</td> <td>equent item sets o much, it uses a special data s ociation informa</td> <td>s without using s a divide-and- tructure named ation.</td> <td></td> <td></td> <td></td>			The FP-Growth Algorithm is candidate generations, thus in conquer strategy. The core of frequent-pattern tree (FP-tree)	an alternative mproving per this method which retains	e way to find fre formance. For s is the usage of a s the item set asso	equent item sets o much, it uses a special data s ociation informa	s without using s a divide-and- tructure named ation.						
• First, it compresses the input database creating an FP-tree instance to represent frequent items. • After this first step, it divides the compressed database into a set of conditional databases, each associated with one frequent pattern. • Finally, each such database is mined separately. Using this strategy, the FP-Growth reduces the search costs by recursively looking for short patterns and then concatenating them into the long frequent patterns			This algorithm works as follows:										
• After this first step, it divides the compressed database into a set of conditional databases, each associated with one frequent pattern. • Finally, each such database is mined separately. Using this strategy, the FP-Growth reduces the search costs by recursively looking for short patterns and then concatenating them into the long frequent patterns			• First, it compresses the input database creating an FP-tree instance to represent frequent items										
$\frac{ \mathbf{x} _{\mathbf{x}} + \mathbf{x} _{\mathbf{x}} + \mathbf{x} _{\mathbf{x}} + \mathbf{x} _{\mathbf{x}} + \mathbf{x} _{\mathbf{x}} + \mathbf{x} _{\mathbf{x}} + \mathbf{x} _{\mathbf{x}} + \mathbf{x} _{\mathbf{x}} + \mathbf{x} _{\mathbf{x}} + \mathbf{x} _{\mathbf{x}} + \mathbf{x} $			• After this first step it divides the compressed database into a set of conditional										
• Finally, each such database is mined separately. Using this strategy, the FP-Growth reduces the search costs by recursively looking for short patterns and then concatenating them into the long frequent patterns $\frac{\hline Transaction D Ems}{\hline T2 (D,E,K,N,O,Y)} \\ \hline T3 (A,E,K,M) \\ \hline T4 (C,K,M,U,Y) \\ \hline T5 (C,E,I,K,O,O) \\ \hline The frequency of each individual item is computed:- \frac{\boxed{Transaction D tems}{\hline T1 (E,K,M,N,O,Y)} \\ \hline T2 (D,E,K,N,O,Y) \\ \hline T3 (A,E,K,M) \\ \hline T4 (C,K,M,U,Y) \\ \hline T5 (C,E,I,K,O,O) \\ \hline N 2 \\ \hline 0 3 \\ \hline U 1 \\ \hline Y 3 \\ \hline \end{pmatrix}$			databases, each associ	ated with one	frequent pattern.								
Using this strategy, the FP-Growth reduces the search costs by recursively looking for short patternsTransaction IDItemsT1{E,K,M,N,O,Y}T2{D,E,K,N,O,Y}T3{A,E,K,M}T4{C,K,M,U,Y}T5{C,E,I,K,O,O}Te frequency of each individual item is computed:-Items CTansaction IDItems T1T1{E,K,M,N,O,Y}T2{D,E,K,N,O,Y}T3{A,E,K,M}T4{C,K,M,U,Y}T5{C,E,I,K,O,O}			• Finally, each such data	abase is mined	l separately								
Transaction ID Items T1 {E,K,M,N,O,Y} T2 {D,E,K,N,O,Y} T3 {A,E,K,M} T4 {C,K,M,U,Y} T5 {C,E,I,K,O,O} Item Frequency A 1 $c_{c,E,I,K,O,O}$ A Item is computed:- D T1 {E,K,M,N,O,Y} T2 {D,E,K,N,O,Y} T3 {A,E,K,M} T4 {C,K,M,U,Y} T5 {C,E,I,K,O,O} M 3 T4 {C,K,M,U,Y} N 2 0 3 U 1 Y 3			Using this strategy, the FP-Gr patterns and then concatenatin	owth reduces t g them into th	the search costs b e long frequent p	by recursively lo patterns	ooking for short						
T1 {E,K,M,N,O,Y} T2 {D,E,K,N,O,Y} T3 {A,E,K,M} T4 {C,K,M,U,Y} T5 {C,E,I,K,O,O} The frequency of each individual item is computed:- Item Frequency T1 {E,K,M,N,O,Y} 1 1 T2 {D,E,K,N,O,Y} I 1 T1 {E,K,M,N,O,Y} K 5 T3 {A,E,K,M} M 3 T4 {C,K,M,U,Y} N 2 T5 {C,E,I,K,O,O} 3 U 1			Transaction	ID		Items							
T2 {D,E,K,N,O,Y} T3 {A,E,K,M} T4 {C,K,M,U,Y} T5 {C,E,I,K,O,O} Item Frequency A 1 item is computed:- A T1 {E,K,M,N,O,Y} T2 {D,E,K,N,O,Y} K 5 T3 {A,E,K,M} T4 {C,K,M,U,Y} N 2 0 3 U 1 V 3			T1		{ <u>E,K</u>	(,M,N,O,Y}							
T3 {A,E,K,M} T4 {C,K,M,U,Y} T5 {C,E,I,K,O,O} The frequency of each individual A item is computed:- A T1 {E,K,M,N,O,Y} T2 {D,E,K,N,O,Y} T3 {A,E,K,M} T4 {C,K,M,U,Y} N 2 0 3 U 1 T4 {C,K,M,U,Y} V 3 U 1 V 3			T2		{ <u>D</u> ,	E,K,N,O,Y}							
T4{C,K,M,U,Y}T5{C,E,I,K,O,O}The frequency of each individual item is computed:- $I tem is computed:-IImage: Transaction IDItemsImage: CT1{E,K,M,N,O,Y}Image: CT2{D,E,K,N,O,Y}Image: C,E,I,K,O,O}T3{A,E,K,M}Image: C,E,I,K,O,O}T5{C,E,I,K,O,O}U1V3U1T5{C,E,I,K,O,O}U1V3$			T3		{/	{ <u>A,E,K,M</u> }							
Ts{C,E,I,K,O,O}The frequency of each individual item is computed:-ItemTransaction IDItemsT1{E,K,M,N,O,Y}T2{D,E,K,N,O,Y}T3{A,E,K,M}T4{C,K,M,U,Y}T5{C,E,I,K,O,O}U1Y3			T4		{ <u>C</u> ,	<u>,K</u> ,M,U,Y}							
The frequency of each individualItemFrequencyitem is computed:- A 1Transaction IDItems D 1T1{E,K,M,N,O,Y} I 1T2{D,E,K,N,O,Y} K 5T3{A,E,K,M} M 3T4{C,K,M,U,Y} N 2T5{C,E,I,K,O,O} O 3U1 Y 3			T5 {C,E,I,K,O,O}										
Interfrequency of each individual A 1 item is computed:- C 2 Transaction ID Items D 1 T1 {E,K,M,N,O,Y} I 1 T2 {D,E,K,N,O,Y} K 5 T3 {A,E,K,M} M 3 T4 {C,K,M,U,Y} N 2 U 1 Y 3			The frequency of each	individual		ltem	Frequency						
item is computed:- C 2 Transaction ID Items D 1 T1 {E,K,M,N,O,Y} I 1 T2 {D,E,K,N,O,Y} K 5 T3 {A,E,K,M} M 3 T4 {C,K,M,U,Y} N 2 U 1 1 Y 3				rinuiviuuai		А	1						
Transaction ID Items T1 {E,K,M,N,O,Y} T2 {D,E,K,N,O,Y} T3 {A,E,K,M} T4 {C,K,M,U,Y} N 2 U 1 V 3 U 1 V 3			item is computed.			C	2						
Transaction ID Items T1 {E,K,M,N,O,Y} T2 {D,E,K,N,O,Y} T3 {A,E,K,M} T4 {C,K,M,U,Y} N 2 O 3 U 1 Y 3			item is computed."			D	1						
T1 {E,K,M,N,O,Y} I 1 T2 {D,E,K,N,O,Y} K 5 T3 {A,E,K,M} M 3 T4 {C,K,M,U,Y} N 2 U 1 U 1 Y 3 Y 3			Transaction ID	1+	ems	E	4						
T2 {D,E,K,N,O,Y} T3 {A,E,K,M} T4 {C,K,M,U,Y} N 2 T5 {C,E,I,K,O,O} U 1 Y 3			T1	{F.K.N	1.N.O.Y}	I	1						
T3 {A,E,K,M} T4 {C,K,M,U,Y} T5 {C,E,I,K,O,O} U 1 Y 3			T2	{D.E.K	(.N.O.Y)	К	5						
T4 {C,K,M,U,Y} N 2 T5 {C,E,I,K,O,O} O 3 U 1 Y 3			T3	<u>{A,E,K,M}</u> { <u>C,K,M,U,Y</u> }		М	3						
T5 {C,E,I,K,O,O} O 3 U 1 Y 3			T4			N	2						
U 1 Y 3			T5 {C,E,I,K,O,O}		0	3							
Y 3					1								
						Y	3						

Item	Frequency
А	1
С	2
D	1
E	4
I	1
К	5
М	3
N	2
0	3
U	1
Y	3

A **Frequent Pattern set (L)** is built which will contain all the elements whose frequency is greater than or equal to the minimum support.

As minimum support be 3.

These elements are stored in descending order of their respective frequencies.

After insertion of the relevant items, the set L looks

like this:- L = {K : 5, E : 4, M : 3, O : 3, Y : 3}

Now, for each transaction, the respective Ordered-Item set is built.

Frequent Pattern set L = {K : 5, E : 4, M : 3, O : 3, Y : 3}

Transaction ID	ltems	Ordered-Item Set
T1	{ <u>E,K</u> ,M <mark>,N</mark> ,O,Y}	{ <u>K,E</u> ,M,O,Y}
T2	{D,E,K,N,O,Y}	{ <u>K,E</u> ,O,Y}
T3	{ <u>A,E</u> ,K,M}	{ <u>K,E</u> ,M}
T4	{C,K,M,U,Y}	{K,M,Y}
T5	{C,E,I,K,O,O}	{K,E,O}

	b)	Give a short example to show that items in a strong association rule actually may be	CO3	L4	7M
		negatively correlated.			
		The support and confidence measures are insufficient at filtering out uninteresting			
		association rules. To tackle this weakness, a correlation measure can be used to augment			
		the support-confidence framework for association rules. This leads to correlation rules of			
		the form			
		$A \rightarrow B$ [support, confidence. correlation].			
		That is, a correlation rule is measured not only by its support and confidence but also by the			
		correlation between itemsets A and B. There are many different correlation measures from			
		which to choose. In this section, we study various correlation measures to determine which			
		would be good for mining large data sets.			
		Lift is a simple correlation measure that is given as follows. The occurrence of itemset A is			
		Independent of the occurrence of itemset B if $D(A \mid P) = D(A)D(P)$			
		$I(A \mid D) = I(A)I(D)$, otherwise itemsets A and B are dependent and correlated as events. This definition can			
		easily be			
		extended to more than two itemsets. The lift between the occurrence of A and B can be			
		measured by computing			
		$lift(A \mid B) = P(A \cup B)/P(A)P(B)$			
		If the resulting value is less than 1, then the occurrence of A is negatively correlated with			
		the occurrence of B. If the resulting value is greater than 1, then A and B are <i>positively</i>			
		<i>correlated</i> , meaning that the occurrence of one implies the occurrence of the other. If the			
		resulting value is equal to 1, then A and B are <i>independent</i> and there is no correlation			
		between them.			
		It assesses the degree to which the occurrence of one "lifts" the occurrence of the other. For			
		example, if A corresponds to the sale of computer games and B corresponds to the sale of			
		videos,			
		then given the current market conditions, the sale of games is said to increase or "lift" the			
		likelihood of the sale of videos by a factor of the value returned by Equation			
8	2)	<u>UIIII-1V</u> Evolain agglomerative and divisive hierarchical clustering algorithm	CO4	12	7M
0	u)	Agglomerative hierarchical clustering. This bottom-up strategy starts by placing each	007		/ 101
		object in its own cluster and then merges these atomic clusters into larger and larger			
		clusters, until all of the objects are in a single cluster or until certain termination conditions			
		are satisfied. Most hierarchical clustering methods belong to this category. They differ only			
		in their definition of intercluster similarity.			
		Divisive hierarchical clustering: This top-down strategy does the reverse of agglomerative			
		hierarchical clustering by starting with all objects in one cluster. It subdivides the cluster			
		into smaller and smaller pieces, until each object forms a cluster on its own or until it			
		satisfies certain termination conditions, such as a desired number of clusters is obtained or			
		the diameter of each cluster is within a certain threshold.			
	b)	Explain about DBSCAN algorithm.	CO4	L2	7M
		DBSCAN (Density-Based Spatial Clustering of Applications with Noise) is a densitybased			
		clustering algorithm. The algorithm grows regions with sufficiently high density into			
		clusters and discovers clusters of arbitrary snape in spatial databases with hoise. It defines a			
		cluster as a maximal set of <i>density-connected</i> points. The basic ideas of density-based			
		and then follow up with an example			
		• The neighborhood within a radius a of a given object is called the a neighborhood of the			
		• The heighborhood within a radius e of a given object is called the e-heighborhood of the object			
		• If the e-neighborhood of an object contains at least a minimum number <i>MinPts</i> of			
		objects, then the object is called a core object.			
		• Given a set of objects, D, we say that an object p is directly density-reachable from			
		object q if p is within the e-neighborhood of q , and q is a core object.			
		• An object p is density-reachable from object q with respect to e and <i>MinPts</i> in a set of			
		objects, D, if there is a chain of objects $p1, :::, pn$, where $p1 = q$ and $pn = p$ such that			
		<i>pi</i> +1 is directly density-reachable from <i>pi</i> with respect to e and <i>MinPts</i> , for 1 _ <i>i</i> _ <i>n</i> , <i>pi</i> 2			
		D.			
		• An object p is density-connected to object q with respect to e and <i>MinPts</i> in a set of			
		objects, D, if there is an object 0 2 D such that both p and q are density-reachable from o with respect to a and MinDes			
		with respect to e and <i>MINTES</i> .			
		relationship is asymmetric. Only core objects are mutually density reachable. Density			
L	1				

		connectivity, however, is a symmetric relat	ion.			
(OR)						
9	a)	Differentiate between density based and gr	id-based methods.	CO4	L2	7M
9	a)	 Density-based methods: Most partitioning methods cluster objects based on the distance between objects. Such methods can find only spherical-shaped clusters and encounter difficulty at discovering clusters of arbitrary shapes. Other clustering methods have been developed based on the notion of <i>density</i>. Their general idea is to continue growing the given cluster as long as the density (number of objects or data points) in the "neighborhood" exceeds some threshold; that is, for each data point within a given cluster the neighborhood of a given radius has to contain atleast a minimum number of points. Such a method can be used to filter out noise (outliers) and discover clusters of arbitrary shape. DBSCAN and its extension, OPTICS, are typical density-based methods that grow clusters objects based on the analysis of the value distributions of density functions. Grid-based methods: Grid-based methods quantize the object space into a finite number of cells that form a grid structure. All of the clustering operations are performed on the grid structure (i.e., on the quantized space). The main advantage of this approach is its fast processing time, which is typically independent of the number of clusters and dependent only on the number of cells in each dimension in the quantized space. STING is a typical example of a grid-based and density-based. 			L2	/1/1
	b)	Both k-means and k-medoids algorithms of the strength and weakness of k-means in co K-means takes the mean of data points to create new points called controids	can perform effective clustering. Illustrate omparison with k-medoids. K-medoids uses points from the data to serve as points called medoids.	CO4	L3	7M
		Centroids. Centroids are new points previously not found in the data. K-means can only by used for numerical data. K-means focuses on reducing the sum of squared distances, also known as the sum of squared error (SSE). K-means is not sensitive to outliers	Medoids are existing points from the data. K-medoids can be used for both numerical and categorical data. K-medoids focuses on reducing the dissimilarities between clusters of data from the dataset. K-medoids is outlier resistant and			
		within the data.K-means is less costly to implement.K-means is faster.	can reduce the effect of outliers. K-medoids is more costly to implement. K-medoids is comparatively not as			
			iast.	1		

Scrutiny Prepared by: P.Ravi Kumar, Dept of IT, BEC HOD,IT