Hal	l Ti	cket	t Nu	mb	er:		

18CSD12

III/IV B.Tech (Regular\Supplementary) DEGREE EXAMINATION

February, 2021		ry, 2021 Computer Science and Enginee	Computer Science and Engineering			
Fift	th Se	Data ware Housing and Data Min	ning			
Tim	e: Th	ree Hours Maximum: 50	Marks			
Ans	wer Q	$Puestion No.1 \ compulsorily. $ (1X10 = 10 N)	Aarks)			
Ansı	wer a	ny one question from each unit. (4X10=40 M	Aarks)			
1.	An	swer all questions (1X10=10 M	Marks)			
	a)	Define Data Warehouse.				
	b)	How to deal with missing values in an attribute?				
	c)	What is the importance of a fact table?				
	a)	What is Data Mining? How to compute confidence measure for an association rule?				
	e) f)	What is Data Mart?				
	g)	Differentiate qualitative and quantitative attributes.				
	h)	Define cluster analysis?				
	i)	What are properties of good Clustering?				
	j)	Define outlier?				
•		UNIT I	7 3 6			
2.	a)	What are Data Mining functionalities? Explain briefly.	5M			
	D)	what is Data Cleaning? Explain various data cleaning tasks.	SIM			
3	a)	Discuss various issues in Data Mining	5M			
5.	b)	Illustrate the Data Transformation by Normalization.	5M			
	,	UNIT II				
4.	a)	Explain Data Warehouse architecture.	5M			
	b)	What is the difference between operational DBMS and Data Warehouse	5M			
-	``	(OR)	7) (
5.	a) b)	Discuss the star and snowflake schema in detail with suitable example.	5M 5M			
	0)		JIVI			
6.	a)	A database has six transactions. Let min-sup = 50% and min-conf = 75%.	5M			
		TID List of items				
		001 Pencil sharpener eraser color papers				
		002 Color papers charts glue sticks				
		002 Color papers, charts, give sticks				
		003 Pencil, glue stick, eraser, pen				
		004 Oil pastels, poster colours, correction tape				
		005 Whitener, pen, pencil, charts, glue stick				
		006 Colour pencils, crayons, eraser, pen				
		Find all frequent item sets using Apriori algorithm. List all the strong association rules.				
	b)	Write the advantages and disadvantages of Apriori and FP-growth Algorithm	5M			
7	-)	(OR)	514			
7.	a) b)	Explain constraint based rule mining	JIVI 5M			
	U)	LINIT IV	JIVI			
8.	a)	What is the goal of clustering? How does partitioning around medoids algorithm achieve this goal?	5M			
0.	b)	Write K-means clustering algorithm.	5M			
		(OR)				
9.		Explain Hierarchical clustering algorithm.	10M			

Hall Ticket Number:

18CSD12

III/IV B.Tech (Regular\Supplementary) DEGREE EXAMINATION

February, 2021

Computer Science and Engineering Data ware Housing and Data Mining Maximum: 50 Marks

Fifth Semester Time: Three Hours

Answer Question No.1 compulsorily.

(1X10 = 10 Marks)

Ansv	ver a	ny one question from each unit. (4X	10=40 Marks
1.	An	iswer all questions (1X	10=10
		Mar	rks)
	a)	Define Data Warehouse.	
		A data warehouse is a subject-oriented, integrated, time-variant, and Non-volatile collection	n of data
		in support of management's decision making process.	
	b)	How to deal with missing values in an attribute?	
		If written any 3 of the below, award 1M	
		1) Ignore the tuple.	
		2) Fill in the missing value manually.	
		3) Use a global constant to fill in the missing value.	
		4) Use the attribute mean to fill in the missing value.	
		5) Use the attribute mean for all samples belonging to the same class as the given tuple.	
		6) Use the most probable value to fill in the missing value	
	c)	What is the importance of a fact table?	
		A fact table consists of two types of columns. The foreign keys column which allows join	ing with
		dimension tables and the measure columns contain the data that is being analyzed.	
	d)	What is Data Mining?	
		Data mining is a process of discovering patterns in large data sets	
	e)	How to compute confidence measure for an association rule?	
		Confidence $(A=>B) = (Number of transactions includes both A and B)/ (Number of transactions)$	sactions
		includes only product A)	
	f)	What is Data Mart?	
		A data mart is a simple form of a data warehouse that is focused on a single subject (or fu	nctional
		area), such as Sales or Finance or Marketing.	
	g)	Differentiate qualitative and quantitative attributes.	
		Quantitative Attributes - Attributes whose values result from counting or measuring something	hing.
		Examples: height, weight, time in the 100 yard dash, number of items sold to a shopper	
		Qualitative Attributes - Attributes that are not measurement variables. Their values do not n	result
		from measuring or counting.	
		Examples: hair color, religion, political party, profession	
	h)	Define cluster analysis?	
		Cluster Analysis means that to find out the group of objects which are similar to each other	in the
		group but are different from the object in other groups.	
	i)	What are properties of good Clustering?	
	,	A good clustering method will produce high quality clusters with	
		• high intra-class similarity	
		• low <u>inter-class</u> similarity	
	i)	Define outlier?	
	5,	An outlier is a data point that differs significantly from other observations.	
		UNIT I	•
2.	a)	What are Data Mining functionalities? Explain briefly.	5M
	Í	Data Mining Functionalities	
		Data mining functionalities are used to specify the kind of patterns to be found in dat	a
		mining tasks.Data mining tasks can be classified into two categories: descriptive and	
		predictive.	
		Descriptive mining tasks characterize the general properties of the data in the databa	se.
		Predictive mining tasks perform inference on the current data in order to make	

predictions.

Concept/Class Description: Characterization and Discrimination

Data can be associated with classes or concepts. For example, in the Electronics store, classes of items for sale include computers and printers, and concepts of customers include bigSpenders and budgetSpenders.

Data characterization

Data characterization is a summarization of the general characteristics or features of a target class of data.

Data discrimination

Data discrimination is a comparison of the general features of target class data objects with the general features of objects from one or a set of contrasting classes.

Mining Frequent Patterns, Associations, and Correlations

Frequent patterns, are patterns that occur frequently in data. There are many kinds of frequent patterns, including itemsets, subsequences, and substructures.

Association analysis

Suppose, as a marketing manager, you would like to determine which items are frequently purchased together within the same transactions.

buys(X,"computer")=buys(X,"software") [support=1%,confidence=50%] where X is a variable representing a customer.Confidence=50% means that if a customer buys a computer, there is a 50% chance that she will buy software as well. Support=1% means that 1% of all of the transactions under analysis showed that computer and software were purchased together.

Classification and Prediction

Classification is the process of finding a model that describes and distinguishes data classes for the purpose of being able to use the model to predict the class of objects whose class label is unknown.

"How is the derived model presented?" The derived model may be represented in various forms, such as classification (IF-THEN) rules, decision trees, mathematical formulae, or neural networks.

A **decision tree** is a flow-chart-like tree structure, where each node denotes a test on an attribute value, each branch represents an outcome of the test, and tree leaves represent classes or class distributions.



Decision tree

Cluster Analysis

In classification and prediction analyze class-labeled data objects, where as clustering analyzes data objects without consulting a known class label.



Cluster Analysis

The objects are grouped based on the principle of maximizing the intraclass similarity and minimizing the interclass similarity. That is, clusters of objects are formed so that objects within a cluster have high similarity in comparison to one another, but are very dissimilar to objects in other clusters.

Outlier Analysis

A database may contain data objects that do not comply with the general behavior or model of the data. These data objects are outliers. Most data mining methods discard

	outliers as noise or exceptions. The analysis of outlier data is referred to as outlier	
	mining.	
b)	What is Data Cleaning? Explain various data cleaning tasks.	5M
	Data Cleaning in Data Mining	
	Quality of your data is critical in getting to final analysis. Any data which tend to be	
	incomplete, noisy and inconsistent can affect your result.	
	Data cleaning in data mining is the process of detecting and removing corrupt or inaccurate	
	records from a record set, table or database.	
	Some data cleaning methods :-	
	1 You can ignore the tuple. This is done when class label is missing. This method is not very	
	effective, unless the tuple contains several attributes with missing values.	
	2 You can fill in the missing value manually. This approach is effective on small data set with	
	some missing values.	
	3 You can replace all missing attribute values with global constant, such as a label like	
	"Unknown" or minus infinity.	
	4 You can use the attribute mean to fill in the missing value. For example customer average	
	Income is 25000 then you can use this value to replace missing value for income.	
	5 Use the most probable value to fill in the missing value.	
	Noisy Data	
	Noise is a random error of variance in a measured variable. Noisy Data may be due to faulty	
	data confection instruments, data entry problems and technology initiation.	
	Rinning:	
	Binning methods sorted data value by consulting its "neighbor- hood" that is the values	
	around it The sorted values are distributed into a number of "buckets" or bins	
	For example	
	Price = 4, 8, 15, 21, 21, 24, 25, 28, 34	
	Partition into (equal-frequency) bins:	
	Bin a: 4. 8. 15	
	Bin b: 21, 21, 24	
	Bin c: 25, 28, 34	
	In this example, the data for price are first sorted and then partitioned into equal-frequency bins	
	of size 3.	
	Smoothing by bin means:	
	Bin a: 9, 9, 9	
	Bin b: 22, 22, 22	
	Bin c: 29, 29, 29	
	In smoothing by bin means, each value in a bin is replaced by the mean value of the bin.	
	Smoothing by bin boundaries:	
	Bin a: 4, 4, 15	
	Bin b: 21, 21, 24	
	Bin c: 25, 25, 34	
	In smoothing by bin boundaries, each bin value is replaced by the closest boundary value.	
	Regression	
	Data can be smoothed by fitting the data into a regression functions.	
	Clustering:	
	Outliers may be detected by clustering, where similar values are organized into groups, or	
	"clusters. Values that fall outside of the set of clusters may be considered outliers.	



- **Mining different kinds of knowledge in databases** Different users may be interested in different kinds of knowledge. Therefore it is necessary for data mining to cover a broad range of knowledge discovery task.
- Interactive mining of knowledge at multiple levels of abstraction The data mining process needs to be interactive because it allows users to focus the search for patterns, providing and refining data mining requests based on the returned results.
- Incorporation of background knowledge To guide discovery process and to express the discovered patterns, the background knowledge can be used. Background knowledge may be used to express the discovered patterns not only in concise terms but at multiple levels of abstraction.
- Data mining query languages and ad hoc data mining Data Mining Query language that allows the user to describe ad hoc mining tasks, should be integrated with a data warehouse query language and optimized for efficient and flexible data mining.
- **Presentation and visualization of data mining results** Once the patterns are discovered it needs to be expressed in high level languages, and visual representations. These representations should be easily understandable.
- Handling noisy or incomplete data The data cleaning methods are required to handle the noise and incomplete objects while mining the data regularities. If the data cleaning methods are not there then the accuracy of the discovered patterns will be poor.

	• Pattern evaluation – The patterns discovered should be interesting because either they represent common knowledge or lack novelty.	
	Performance Issues	
	There can be performance-related issues such as follows –	
	• Efficiency and scalability of data mining algorithms – In order to effectively extract the information from huge amount of data in databases, data mining algorithm must be efficient and scalable	
	 Parallel, distributed, and incremental mining algorithms – The factors such as huge size of databases, wide distribution of data, and complexity of data mining methods motivate the development of parallel and distributed data mining algorithms. These algorithms divide the data into partitions which is further processed in a parallel fashion. Then the results from the partitions is merged. The incremental algorithms, update databases without mining the data again from scratch. 	
	 Diverse Data Types Issues Handling of relational and complex types of data – The database may contain complex data objects, multimedia data objects, spatial data, temporal data etc. It is not possible for one system to mine all these kind of data. Mining information from heterogeneous databases and global information systems – The data is available at different data sources on LAN or WAN. These data sources may be structured, semi structured or unstructured. Therefore mining the 	
	knowledge from them adds challenges to data mining.	
b)	Illustrate the Data Transformation by Normalization. The data transformation is a basic element of data mining. It means transforming the data, namely converting the source data in to another format that allows processing data effectively. The main purpose of data normalization is to minimize or even exclude duplicated data. This is a very essential and important issue because it is increasingly problematic to keep in data in relational databases, which store identical data in more than one place.	5M
	Min-Max normalization The first technique we will cover is min-max normalization. It is the linear transformation of the original unstructured data. It scales the data from 0 to 1. It is calculated by the following formula:	
	$v' = rac{v - \min F}{\max F - \min F} (new_max_F - new_min_F) + new_min_F$,	
	where is the current value of feature <i>F</i> . Let us consider one example to make the calculation method clear. Assume that the minimum and maximum values for the feature F are \$50,000 and \$100,000 correspondingly. It needs to range <i>F</i> from 0 to 1. In accordance with min-max normalization, $v = $80,000$ is transformed to: $v' = \frac{80,000 - 50,000}{100,000 - 50,000} + (1 - 0) + 0 = \frac{3}{5} = 0,6$	
	As you can see this technique enables to interpret the data easily. There are no large numbers, only concise data that do not require further transformation and can be used in decision-making process immediately. Z-score normalization The next technique is z-score normalization. It is also called zero-mean normalization. The	
	essence of this technique is the data transformation by the values conversation to a common scale where an average number equals zero and a standard deviation is one. A value is normalized to ' under the formula:	
	$v' = \frac{v - \overline{F}}{\sigma_F},$	
	Here is the mean and is the standard deviation of feature F	
	Here is an example of the calculation of a value.	
	On the supposition that the mean of feature is \$65,000 and its standard deviation is \$ 18,000. Applying the z-score normalization we get the following mean of the value equals to \$85,800:	



		 repository, which stores information about 2. The middle tier is an OLAP server that OLAP (ROLAP) model, that is, an emultidimensional data to standard relat (MOLAP) model, that is, a special-purp data and operations 3. The top tier is a front-end client layer tools, and/or data mining tools (e.g., trender) 	It the data warehouse and its contents. At is typically implemented using either (1 extended relational DBMS that maps op tional operations; or (2) a multidimens pose server that directly implements mult er, which contains query and reporting to a analysis, prediction, and so on).) a relational perations on ional OLAP idimensional pols, analysis	
	b)	What is the difference between operational	al DBMS and Data Warehouse	5	5M
		Operational Database Systems	Data Warehouses		
		Operational systems are generally designed to support high-volume transaction processing.	Data warehousing systems are generally designed to support high-volume analytical processing. (i.e. OLAP).		
		Operational systems focuses on Data in.	Data warehousing systems focuses on Information out.		
		In Operational systems data is stored with a functional or process orientation.	In Data warehousing systems data is stored with a subject orientation.		
		Performance is low for analysis queries.	Performance is high for analysis queries.		
		It is used for Online Transactional Processing (OLTP)	It is used for Online Analytical Processing (OLAP).		
		Operational systems represent current transactions.	Data warehousing systems reads the historical data.		
		Data within operational systems are generally updated regularly.	Data within a data warehouse is non-volatile, meaning when new data is added old data is not erased so rarely updates.		
		Complex data structures.	Multi dimensional data structures.		
5	-)	Disconstitution of the second	(OR)		514
5.	a)	Star Schema: Star schema is the type of multidimensio schema, The fact tables and the dimensio key join is used. This schema forms a sta	nal model which is used for data warehous on tables are contained. In this schema few or with fact table and dimension tables.	se. In star ver foreign-	51 v1
		Dimension Table	Dimension Table		
		Fact Table			
		Table	Dimension Table		
		Dimension Table			
		Snowflake Schema:			





	r –		1
		$\label{eq:second} \begin{split} \end{tabular} & tab$	
	b)	Write the advantages and disadvantages of Apriori and FP-growth Algorithm	5M
		Apriori Algorithm	
		Easy to understand algorithm	
		• Join and Prune steps are easy to implement on large itemsets in large databases	
		 Disadvantages: It requires high computation if the itemsets are very large and the minimum support is 	
		kept very low.	
		• The entire database needs to be scanned. FP-Growth Algorithm	
		Advantages:	
		• Faster than Apriori algorithm.	
		No candidate generation.	
		• Unly two passes over dataset.	
		• FP tree may not fit in memory.	
		• FP tree is expensive to build.	
7.	a)	(OK) What are different methods to improve Apriori algorithm's efficiency?	5M
		• Hash-based itemset counting: A k -itemset whose corresponding hashing bucket count is	
		 Transaction reduction: A transaction that does not contain any frequent k-itemset is 	
		useless in subsequent scans	
		• Partitioning: Any itemset that is potentially frequent in DB must be frequent in at least one of the partitions of DB	
		• Sampling: mining on a subset of given data, lower support threshold + a method to	

		determine the completeness	
		• Dynamic itemset counting: add new candidate itemsets only when all of their subsets are	
		estimated to be frequent	
	b)	Explain constraint based rule mining.	5M
		• A data mining process may uncover thousands of rules from a given set of data, most of	
		which end up being unrelated or uninteresting to the users.	
		• Often, users have a good sense of which "direction" of mining may lead to interesting	
		patterns and the "form" of the patterns or rules they would like to find.	
		• Thus, a good heuristic is to have the users specify such intuition or expectations as	
		Constraints to confine the search space.	
		Constraint based mining provides	
		• Constraint based mining provides	
		• System Ontimization: explores constraints to help efficient mining	
		• The constraints can include the following:	
		 Knowledge type constraints: These specify the type of knowledge to be mined, such as 	
		association or correlation	
		• Data constraints . These specify the set of task-relevant data	
		• Dimension/level constraints: These specify the desired dimensions (or	
		attributes) of the data, or levels of the concept hierarchies, to be used in mining.	
		• Interestingness constraints: These specify thresholds on statistical measures of rule	
		interestingness, such as support, confidence, and correlation.	
		• Rule constraints: These specify the form of rules to be mined. Such constraints	
		may be expressed as rule templates, as the maximum or minimum number of	
		predicates that can occur in the rule antecedent or consequent, or as relationships	
		among attributes, attribute values, and/or aggregates. The above constraints can	
		be specified using a high-level declarative data mining query language and user	
		interface.	
		Constraint based association rules: - In order to make the mining process more efficient rule	
		based constraint mining allows users to describe the rules that they would like to uncover	
		specified by the user encourages interactive exploratory mining and analysis	
		specified by the user encourages interactive exploratory mining and anarysis.	
	1	UNIT IV	
8.	a)	What is the goal of clustering? How does partitioning around medoids algorithm achieve this	5M
	, i i i i i i i i i i i i i i i i i i i	goal?	
		The goal of clustering is to identify distinct groups in a dataset.	
		PAM algorithm identifies distinct groups in a dataset using mediods.	
		 It is used Find representative objects, called <u>medoids</u>, in clusters 	
		• <i>PAM</i> (Partitioning Around Medoids)	
		 starts from an initial set of medoids and iteratively replaces one of the medoids 	
		by one of the non-medoids if it improves the total distance of the resulting	
		clustering	
		- PAM works effectively for small data sets, but does not scale well for large data	
		Sets	
1		CLARA CLAPANS · Dandomized sampling	
		• CLARANS :Randonnized sampling The above said are the three algorithms which are used in DAM	
		Focusing \pm spatial data structure (Ester et al. 1995)	
		Tocusing + spanar data structure (Ester et al., 1995)	
<u> </u>	b)	Write K-means clustering algorithm	5M
1	- /	Given k, the k-means algorithm is implemented in 4 steps:	
		- Partition objects into k nonempty subsets	
		- Compute seed points as the centroids of the clusters of the current partition. The	
		centroid is the center (mean point) of the cluster.	
		 Assign each object to the cluster with the nearest seed point. 	
		- Go back to Step 2, stop when no more new assignment.	
		- Go back to Step 2, stop when no more new assignment. With relevant example	
		- Go back to Step 2, stop when no more new assignment. With relevant example	
0		- Go back to Step 2, stop when no more new assignment. With relevant example (OR)	101/

