

(Duration : 3 hrs)

- Note: 1) Question No. 1 is **compulsory**  
 2) Attempt any **Four** questions from the remaining **Six** questions.

Q 1.	a)	<p>Briefly compare <b>ANY TWO</b> following concepts                  a)DW and Data mart</p> <table border="1" data-bbox="256 449 1430 1287"> <thead> <tr> <th data-bbox="256 449 841 510"><b>Datamart</b></th> <th data-bbox="841 449 1430 510"><b>Data Warehouse</b></th> </tr> </thead> <tbody> <tr> <td data-bbox="256 510 841 730">Data mart is usually sponsored at the department level and developed with a specific issue or subject in mind, a data mart is a data warehouse with a focused objective.</td> <td data-bbox="841 510 1430 730">Data warehouse is a “Subject-Oriented, Integrated, Time-Variant, Nonvolatile collection of data in support of decision making”.</td> </tr> <tr> <td data-bbox="256 730 841 831">A data mart is used on a business division/ department level.</td> <td data-bbox="841 730 1430 831">A data warehouse is used on an enterprise level.</td> </tr> <tr> <td data-bbox="256 831 841 1047">A Data Mart is a subset of data from a Data Warehouse. Data Marts are built for specific user groups.</td> <td data-bbox="841 831 1430 1047">A Data Warehouse is simply an integrated consolidation of data from a variety of sources that is specially designed to support strategic and tactical decision making.</td> </tr> <tr> <td data-bbox="256 1047 841 1224">By providing decision makers with only a subset of data from the Data Warehouse, Privacy, Performance and Clarity Objectives can be attained.</td> <td data-bbox="841 1047 1430 1224">The main objective of Data Warehouse is to provide an integrated environment and coherent picture of the business at a point in t</td> </tr> <tr> <td data-bbox="256 1224 841 1287">Uses bottom-up approach</td> <td data-bbox="841 1224 1430 1287">Uses top-down approach</td> </tr> </tbody> </table>	<b>Datamart</b>	<b>Data Warehouse</b>	Data mart is usually sponsored at the department level and developed with a specific issue or subject in mind, a data mart is a data warehouse with a focused objective.	Data warehouse is a “Subject-Oriented, Integrated, Time-Variant, Nonvolatile collection of data in support of decision making”.	A data mart is used on a business division/ department level.	A data warehouse is used on an enterprise level.	A Data Mart is a subset of data from a Data Warehouse. Data Marts are built for specific user groups.	A Data Warehouse is simply an integrated consolidation of data from a variety of sources that is specially designed to support strategic and tactical decision making.	By providing decision makers with only a subset of data from the Data Warehouse, Privacy, Performance and Clarity Objectives can be attained.	The main objective of Data Warehouse is to provide an integrated environment and coherent picture of the business at a point in t	Uses bottom-up approach	Uses top-down approach	(10 )  (05 marks each)  (1mark for each point)
<b>Datamart</b>	<b>Data Warehouse</b>														
Data mart is usually sponsored at the department level and developed with a specific issue or subject in mind, a data mart is a data warehouse with a focused objective.	Data warehouse is a “Subject-Oriented, Integrated, Time-Variant, Nonvolatile collection of data in support of decision making”.														
A data mart is used on a business division/ department level.	A data warehouse is used on an enterprise level.														
A Data Mart is a subset of data from a Data Warehouse. Data Marts are built for specific user groups.	A Data Warehouse is simply an integrated consolidation of data from a variety of sources that is specially designed to support strategic and tactical decision making.														
By providing decision makers with only a subset of data from the Data Warehouse, Privacy, Performance and Clarity Objectives can be attained.	The main objective of Data Warehouse is to provide an integrated environment and coherent picture of the business at a point in t														
Uses bottom-up approach	Uses top-down approach														
b) Semi Joins and Bloom Joins		<table border="1" data-bbox="256 1350 1430 1978"> <thead> <tr> <th data-bbox="256 1350 841 1392"><b>Semi join</b></th> <th data-bbox="841 1350 1430 1392"><b>bloom join</b></th> </tr> </thead> <tbody> <tr> <td data-bbox="256 1392 841 1539">Semi join is a method used for efficient query processing in a distributed database environments.</td> <td data-bbox="841 1392 1430 1539">bloom join is another method used to avoid transferring unnecessary data between sites when executing queries in a distributed database environments</td> </tr> <tr> <td data-bbox="256 1539 841 1686">Only some of the attributes (or tuples) that are required for the join need to be transferred between the sites to execute the query efficiently</td> <td data-bbox="841 1539 1430 1686">rather than transferring the join column itself, a compact representation of the join column is transferred between the sites</td> </tr> <tr> <td data-bbox="256 1686 841 1791">Semi join is a method that can be used to reduce the amount of data shipped between the sites.</td> <td data-bbox="841 1686 1430 1791">Bloom join is a method that can be used to reduce the amount of data shipped between the sites.</td> </tr> <tr> <td data-bbox="256 1791 841 1978">In semi join, only the join column is transferred from one site to the other and then that transferred column is used to reduce the size of the shipped relations between the other sites.</td> <td data-bbox="841 1791 1430 1978">Bloom join uses a bloom filter which employs a bit vector to execute membership queries.</td> </tr> </tbody> </table>	<b>Semi join</b>	<b>bloom join</b>	Semi join is a method used for efficient query processing in a distributed database environments.	bloom join is another method used to avoid transferring unnecessary data between sites when executing queries in a distributed database environments	Only some of the attributes (or tuples) that are required for the join need to be transferred between the sites to execute the query efficiently	rather than transferring the join column itself, a compact representation of the join column is transferred between the sites	Semi join is a method that can be used to reduce the amount of data shipped between the sites.	Bloom join is a method that can be used to reduce the amount of data shipped between the sites.	In semi join, only the join column is transferred from one site to the other and then that transferred column is used to reduce the size of the shipped relations between the other sites.	Bloom join uses a bloom filter which employs a bit vector to execute membership queries.	(05 marks each)  (1mark for each point)		
<b>Semi join</b>	<b>bloom join</b>														
Semi join is a method used for efficient query processing in a distributed database environments.	bloom join is another method used to avoid transferring unnecessary data between sites when executing queries in a distributed database environments														
Only some of the attributes (or tuples) that are required for the join need to be transferred between the sites to execute the query efficiently	rather than transferring the join column itself, a compact representation of the join column is transferred between the sites														
Semi join is a method that can be used to reduce the amount of data shipped between the sites.	Bloom join is a method that can be used to reduce the amount of data shipped between the sites.														
In semi join, only the join column is transferred from one site to the other and then that transferred column is used to reduce the size of the shipped relations between the other sites.	Bloom join uses a bloom filter which employs a bit vector to execute membership queries.														

		<p>No such filters are used in semijoin</p>	<p>a bloom filter is built using the join column and it is transferred between the sites and then the joining operations are performed.</p>	<p>(05 marks each)</p> <p>(1mark for each point)</p>
<p>c) OODBMS and ORDBMS</p>				
<p><b>OODBMS</b></p>		<p><b>ORDBMS</b></p>		
<p>Supports object oriented features extensively</p>	<p>Limited Support to object oriented</p>			
<p>Relatively less performance</p>	<p>Expected to perform well</p>			
<p>Query language is ODL/OQL</p>	<p>Query language is SQL3</p>			
<p>OODBMSs try to add DBMS functionality to a programming language</p>	<p>ORDBMSs try to add richer data types to a relational DBMS</p>			
<p>OODBMS put more emphasis on the role of the client side This can improve long, process intensive, transactions</p>	<p>ORDBMS SQL is still the language for data definition, manipulation and query.</p>			
<p>OODBMS have been optimized to directly support object-oriented applications and specific OO languages</p>	<p>ORDBMS Most third-party database tools are written for the relational model and will therefore be compatible with SQL3</p>			
<p>We can use OODBMS In applications that generally retrieve relatively few (generally physically large) highly complex objects and work on them for long periods of time.</p>	<p>We can use ORDBMS In applications that process a large number of short-lived (generally ad-hoc query) transactions on data items that can be complex in structure</p>			
<p>b)</p>	<p>Write Short Notes on <b>ANY TWO</b></p> <p>a)Web Mining <span style="float: right;">(05 marks)</span></p> <ul style="list-style-type: none"> <li>➤ Definition of web mining</li> <li>➤ Types of web mining <ul style="list-style-type: none"> <li>○Content mining</li> <li>○Structure mining</li> <li>○Usage mining</li> </ul> </li> </ul> <p>b)Neural networks <span style="float: right;">(05 marks)</span></p> <ul style="list-style-type: none"> <li>➤ Definition of neural network</li> <li>➤ Structure of neural network (diagram)</li> <li>➤ Steps in neural network</li> <li>➤ Backpropagation</li> </ul> <p>c)Spatial and Geographic DB <span style="float: right;">(05 marks)</span></p> <ul style="list-style-type: none"> <li>➤ A <b>spatial database</b> is a database that is optimized to store and query data that represents objects defined in a geometric space. Most spatial databases allow representing simple geometric objects such as points, lines and polygons. Some spatial databases handle more complex structures such as 3D objects, topological coverages, linear networks.</li> <li>➤ A geodatabase (also <b>geographical database</b> and <b>geospatial database</b>) is a <b>database</b> of <b>geographic</b> data, such as countries, administrative divisions, cities, and related information. Such <b>databases</b> can be useful for websites that wish to identify the locations of their visitors for customization purposes.</li> <li>➤ Features of spatial database</li> <li>➤ Spatial database elements</li> <li>➤ Application of geographic data</li> </ul>			<p>(10)</p>

Q 2.	a)	<p>Explain concurrency control and recovery process in Distributed Database Management System. (08)</p> <p><b>DISTRIBUTED CONCURRENCY CONTROL (4 marks)</b></p> <p>Ans:</p> <ul style="list-style-type: none"> <li>➤ Lock management can be distributed across sites in many ways:</li> <li>➤ Centralized: A single site is in charge of handling lock and unlock requests for all objects.</li> <li>➤ Primary copy: One copy of each object is designated as the primary copy. All requests to lock or unlock a copy of this object are handled by the lock manager at the site where the primary copy is stored, regardless of where the copy itself is stored.</li> <li>➤ Fully distributed: Requests to lock or unlock a copy of an object stored at a site are handled by the lock manager at the site where the copy is stored.</li> </ul> <p><b>DISTRIBUTED RECOVERY (4 marks)</b></p> <ul style="list-style-type: none"> <li>➤ Two-Phase Commit (2PC) protocol</li> <li>➤ Three-Phase Commit</li> </ul>	(08)
	b)	<p>Explain bitmap index with example. When does it make sense to use bitmap index? (07)</p> <p>Ans:</p> <p><b>What is bitmap index (3 marks)</b></p> <ul style="list-style-type: none"> <li>➤ <b>bitmap index</b> is a special kind of database index that uses bitmaps.</li> <li>➤ Bitmap indexes have traditionally been considered to work well for <i>low-cardinality columns</i>, which have a modest number of distinct values, either absolutely, or relative to the number of records that contain the data.</li> <li>➤ The extreme case of low cardinality is Boolean data which has two values, True and False.</li> <li>➤ Bitmap indexes use bit arrays (commonly called bitmaps) and answer queries by performing bitwise logical operations on these bitmaps.</li> <li>➤ Bitmap indexes have a significant space and performance advantage over other structures for query of such data.</li> <li>➤ Their drawback is they are less efficient than the traditional B tree indexes for columns whose data is frequently updated</li> </ul> <p><b>Example of bitmap index (2 marks)</b></p> <p><b>When does it make sense to use a bitmap index? (2 marks)</b></p> <p>Bitmap indexes are meant to be used on low cardinality columns. A low cardinality column just means that the column has relatively few unique values. For example, a column called Sex which has only “Male” and “Female” as the two possible values is considered low cardinality because there are only two unique values in the column.</p>	(07)
Q 3.	a)	<p>What is Data Warehouse? Why it is needed? Explain ETL process in data warehousing. (08)</p> <p>Ans:</p> <p>Data Warehouse Definition. (1 marks)</p> <p>Use of data warehouse (2 marks)</p> <p>ETL process in brief (5 marks)</p> <ul style="list-style-type: none"> <li>➤ Full form of ETL</li> <li>➤ Extraction: various sources</li> <li>➤ Transformation : Techniques</li> <li>➤ Loading : Types of loading</li> </ul>	(08)

	b)	What are frequent itemsets? Describe an algorithm for finding frequent itemsets.	(07)										
	Ans:	Frequent itemsets definition (2 marks) Example for frequent itemset (1 marks) Apriori algorithm (4 marks)											
Q 4.	a)	Define Clustering in data mining. Explain K-Mean clustering with suitable example.	(08)										
	Ans:	Defination of clustering (2 marks) Algorithm for K-mean clustering (4 marks) Example (2 marks)											
	b)	Explain ORDBMS Implementation challenges in detail.	(07)										
	Ans:	ORDBMS definition (1 marks) Implementation challenges are as follows: (6 marks) 1) <u>Storage and Access Methods</u> 2) <u>Indexing New Types:</u> 3) <u>QUERY PROCESSING:</u> 4) <u>User-defined aggregation functions:</u> 5) <u>Method Security</u> 6) <u>Method catching</u>											
Q 5.	a)	Define OLAP. Explain MOLAP and ROLAP system with suitable diagram.	(08)										
	Ans:	Definition of OLAP (online Analytical Processing) Definition of MOLAP(Multi-dimensional Online Analytical Processing) MOLAP model with suitable diagram (3 marks) Definition of ROLAP(Multi-dimensional Online Analytical Processing) ROLAP model with suitable diagram (3 marks) Characteristics of MOLAP and ROLAP (2 marks)											
	b)	Explain features of XML. Differentiate between XML and HTML.	(07)										
	Ans:	XML definition and features (2 marks) Features: 1)self descriptive 2)free and extensible 3)platform independent 4)provides domain specific vocabulary 5)it allows data interchange 6)smart searches  Difference between HTML and XML (5 marks)											
		<table border="1"> <thead> <tr> <th>HTML</th> <th>XML</th> </tr> </thead> <tbody> <tr> <td>HTML is an abbreviation for HyperText Markup Language.</td> <td>XML stands for eXtensible Markup Language.</td> </tr> <tr> <td>HTML was designed to display data with focus on how data looks.</td> <td>XML was designed to be a software and hardware independent tool used to transport and store data, with focus on what data is.</td> </tr> <tr> <td>HTML is a markup language itself.</td> <td>XML provides a framework for defining markup languages.</td> </tr> <tr> <td>HTML is a presentation language.</td> <td>XML is neither a programming language nor a presentation language.</td> </tr> </tbody> </table>	HTML	XML	HTML is an abbreviation for HyperText Markup Language.	XML stands for eXtensible Markup Language.	HTML was designed to display data with focus on how data looks.	XML was designed to be a software and hardware independent tool used to transport and store data, with focus on what data is.	HTML is a markup language itself.	XML provides a framework for defining markup languages.	HTML is a presentation language.	XML is neither a programming language nor a presentation language.	
HTML	XML												
HTML is an abbreviation for HyperText Markup Language.	XML stands for eXtensible Markup Language.												
HTML was designed to display data with focus on how data looks.	XML was designed to be a software and hardware independent tool used to transport and store data, with focus on what data is.												
HTML is a markup language itself.	XML provides a framework for defining markup languages.												
HTML is a presentation language.	XML is neither a programming language nor a presentation language.												

		HTML is not case sensitive.	XML is case sensitive.	
		HTML is used for designing a web-page to be rendered on the client side.	XML is used basically to transport data between the application and the database.	
		HTML has its own predefined tags.	While what makes XML flexible is that custom tags can be defined and the tags are invented by the author of the XML document.	
		HTML is not strict if the user does not use the closing tags.	XML makes it mandatory for the user to close each tag that has been used.	
		HTML does not preserve white space.	XML preserves white space.	
		HTML is about displaying data, hence static.	XML is about carrying information, hence dynamic.	

Q 6.

a)

Find out the association rules with support 50% and confidence 70% from the following sample data.

(08)

Transactions	Items
T1	Bread, Jelly, Milk
T2	Butter, Jelly, Juice
T3	Bread, Butter, Jelly, Juice
T4	Butter, Juice

**Transactions**

Transactions	Items
T <sub>1</sub>	Bread, Jelly, milk
T <sub>2</sub>	Butter, Jelly, Juice
T <sub>3</sub>	Bread, Butter, Jelly, Juice
T <sub>4</sub>	Butter, Juice

**Solution:**

Step 1: Scan D for count of each candidate.  
The candidate list is  
{ Bread, Butter, Jelly, milk, Juice }

C <sub>1</sub> =	Itemset	sup
	{ Bread }	2
	{ Butter }	3
	{ Jelly }	3
	{ milk }	1
	{ Juice }	3

Step 2: Compare candidate support count with min\_support count (ie 50%)

L <sub>1</sub> =	Itemset	sup
	{ Bread }	2
	{ Butter }	3
	{ Jelly }	3
	{ Juice }	3

step 3: Generate candidate  $c_2$  from  $L_1$

$c_2 =$	Itemset
	{ Bread, Butter }
	{ Bread, Jelly }
	{ Bread, Juice }
	{ Butter, Jelly }
	{ Butter, Juice }
	{ Jelly, Juice }

step 4: Scan  $D$  for count of each candidate in  $c_2$  and find the support.

$c_2 =$	Itemset	SUP
	{ Bread, Butter }	1
	{ Bread, Jelly }	2
	{ Bread, Juice }	1
	{ Butter, Jelly }	2
	{ Butter, Juice }	3
	{ Jelly, Juice }	2

step 5: Compare candidate  $c_2$  support count with the min-support count.

$L_2 =$	Itemset	SUP
	{ Bread, Jelly }	2
	{ Butter, Jelly }	2
	{ Butter, Juice }	3
	{ Jelly, Juice }	2

step 6: Generate candidate c3 from L2

c3 =	Itemset
	{ Bread, Jelly, Juice }
	{ Butter, Jelly, Juice }
	{ Bread, Butter, Jelly }

step 7: Scan D for count of each candidate in c3.

	Itemset	Sup
c3 =	{ Bread, Jelly, Juice }	1
	{ Butter, Jelly, Juice }	2
	{ Bread, Butter, Jelly }	1

Step 8: Compare candidate c3 support count with the min-support count.

	Itemset	sup
L3 =	{ Butter, Jelly, Juice }	2

step 9: SO data contain the frequent itemset L { Butter, Jelly, Juice }

Therefore the association rule that can be generated from L3 are as follows:

$$\text{confidence } (A \Rightarrow B) = \frac{A \cup B}{A}$$

Association Rule	Sup	Confidence	confidence%
Butter $\wedge$ Jelly $\Rightarrow$ Juice	2	$2/2 = 1$	100%
Jelly $\wedge$ Juice $\Rightarrow$ Butter	2	$2/2 = 1$	100%
Butter $\wedge$ Juice $\Rightarrow$ Jelly	2	$2/3 = 0.66$	66%
Butter $\Rightarrow$ Jelly $\wedge$ Juice	2	$2/3 = 0.66$	66%
Jelly $\Rightarrow$ Butter $\wedge$ Juice	2	$2/3 = 0.66$	66%
Butter $\Rightarrow$ Juice $\Rightarrow$ Butter $\wedge$ Jelly	2	$2/3 = 0.66$	66%

The minimum confidence threshold is 70%.

$\therefore$  Final rules are

Rule 1: Butter  $\wedge$  Jelly  $\Rightarrow$  Juice

Rule 2: Jelly  $\wedge$  Juice  $\Rightarrow$  Butter

b) Define the term fragmentation and replication in terms of where data is stored and also how the objects are uniquely identified in distributed database.

(07)

Ans: Definition of fragmentation and replication **(3 marks)**  
Types of fragmentation  
➤ Horizontal and vertical  
Distributed catalog management **(2 marks)**  
➤ Naming object  
➤ Local name field  
➤ Birth site field  
➤ Global name field  
Distributed data independence **(2 marks)**

Q 7. a) What are different complex data types available in ORDBMS? Explain with example.

(08)

Ans: There are many complex data types available in ORDBMS:  
Some of them are listed below-- (any four can be explained) **each carry (2 marks)** with suitable example and syntax.  
➤ Collection data types: LIST, SET, MULTISSET  
➤ Abstract data type (ADT)  
➤ User defined types (UDT)  
➤ References (REF and Deref), OID  
➤ ROW type

		➤ VARRAY and NESTED TABLE	
	b)	Explain Star schema, snowflake schema and fact constellation with suitable example.	(07)
	Ans:	<p>Star schema <b>(2 marks)</b></p> <ul style="list-style-type: none"> <li>➤ Structure of star schema</li> <li>➤ Example</li> </ul> <p>snowflake schema <b>(2 marks)</b></p> <ul style="list-style-type: none"> <li>➤ Structure of snowflake schema</li> <li>➤ Example</li> </ul> <p>Advantages and disadvantages of star and snowflake schema <b>(1 marks)</b></p> <p>fact constellation <b>(2 marks)</b></p> <ul style="list-style-type: none"> <li>➤ Structure of snowflake schema</li> <li>➤ Example</li> </ul>	