

Data Warehousing and Business Intelligence (DS3003)

Date: September 23rd 2024

Course Instructor(s)

M. Ishaq Raza

Sessional-I Exam

Total Time (Hrs.): 1
Total Marks: 25
Total Questions: 3

Roll No

Section

Student Signature

Solution

Do not write below this line.

Note: Please ensure that you attempt all questions and their respective parts in the given order.

Consider the following case study for the next two questions:

Bill Date Dim: Bill Date, Bill Day Desc, Bill Month ID, Bill Month Desc, Bill Year ID, Bill Year Desc

Customer Dim: Customer Code, Customer Desc, City ID, City Desc, Country ID, Country Desc

Sales Rep Dim: Sales Rep No, Sales Rep Desc, Channel ID, Channel Desc

Rate Plan Dim: Rate Plan Id, Rate Plan Desc, Rate Plan Type Code, Rate Plan Type Desc

Billing Fact: Bill Date, Customer Code, Sales Rep No, Rate Plan Id, No of Calls, No of Total Minutes, Taxes, Regulatory Charge

Assume: 10,000 customers, 200 cities, 5 countries, 40 sales rep, 4 channels, 30 rate plans, 3 rate plan types, and 3 years billing history.

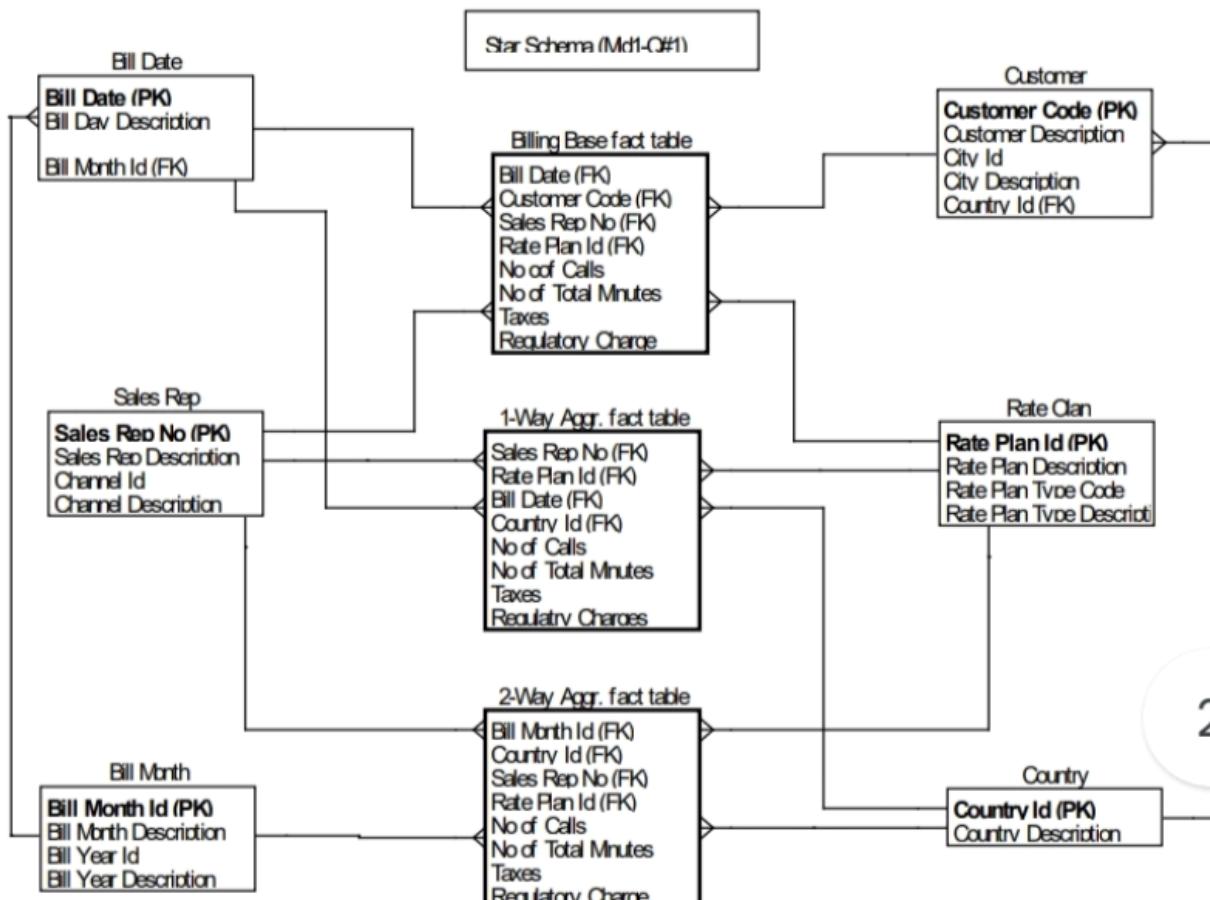
National University of Computer and Emerging Sciences

Lahore Campus

CLO # 2: Demonstrate an understanding of the fundamental concepts of the Star and the Snowflake Schema; learn how to design the schema of a DW based on these two models.

Q. No 1: Draw the appropriate star schema that includes a base fact table, a one-way aggregate fact table (along Customer Country), and a two-way aggregate fact table (along Bill Month and Customer Country). Show the primary keys, foreign keys and all the relationships between the dimensions and fact tables. Note: Draw only one diagram that includes base fact table as well as aggregate fact tables. [10]

Ans:



2/4

National University of Computer and Emerging Sciences
Lahore Campus

CLO # 2: Demonstrate an understanding of the fundamental concepts of the Star and the Snowflake Schema; learn how to design the schema of a DW based on these two models.

Q. No 2: Estimate the size (in number of rows) of the above customer dimension table, sales rep dimension table, billing base fact table, and both the aggregate fact tables. [5]

Ans:

Customer Dimension: $10000+5= 10005$ rows

Sales Rep Dimension: 40 rows

Base Fact Table: $(3 \times 365) \times 10000 \times 40 \times 30 = 13,140,000,000$ rows

Aggregate Fact Table1 (along country):

3×365 (day) $\times 5$ (country) $\times 40$ (sales rep) $\times 30$ (rate plan) = 6,570,000 rows

Aggregate Fact Table2 (along month & country):

3×12 (month) $\times 5$ (country) $\times 40$ (sales rep) $\times 30$ (rate plan) = 216,000 rows

3/4

**National University of Computer and Emerging Sciences
Lahore Campus**

CLO # 2: Demonstrate an understanding of the fundamental concepts of the Star and the Snowflake Schema; learn how to design the schema of a DW based on these two models.

Q. No 3: Briefly answer the following questions. [10]

- a. Pick any one architecture for building a data warehouse and list the advantages and disadvantages of that architecture.
- b. What are the different types of OLAP? Which type of OLAP can handle large amounts of data? Justify your answer.
- c. How does a snowflake schema differ from a star schema? Name two advantages of the snowflake schema.
- d. Differentiate between pre-join denormalization and column-replication denormalization techniques. Explain with an example.
- e. When would you use partitioned cubes in multidimensional online analytical processing (MOLAP)?

Ans: See Textbook/Notes.

Q1)	Bill Date Dim:	Customer Dim:	Sales Rep Dim:
	<u>Bill Date</u>	<u>Customer Code</u>	<u>Sales Rep No</u>
	<u>Bill Day Desc</u>	<u>Customer Desc</u>	<u>Sales Rep Desc</u>
	<u>Bill Month ID</u>	<u>City ID</u>	<u>Channel ID</u>
	<u>Bill Month Desc</u>	<u>City Desc</u>	<u>Channel Desc</u>
	<u>Bill Year ID</u>	<u>Country ID</u>	
	<u>Bill Year Desc</u>	<u>Country Desc</u>	

Rate Plan Dim:
RatePlanID
RatePlanDesc
RatePlanTypeCode
RatePlanTypeDesc

Billing Fact:

BillDate
CustomerCode
SalesRepNo
RatePlanID
NoOfCalls
No. of total Minutes
Taxes
RegulatoryChange

→ 10,000 customers

200 cities

5 countries

40 Sales rep

4 channels

30 rate plans

3 rate plan types

3 years billing history

1 year 12

(Q2)

5

Customer Dim: 10,000 rows

Sales Rep Dim: 40 rows

Billing Base Fact: $3 \times 365 \times 10,000 \times 40 \times 30$
 $: 1.314 \times 10^{10}$ rows

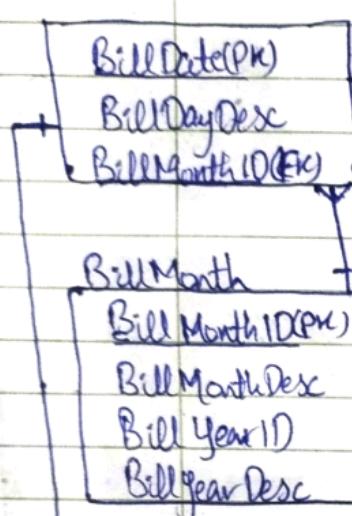
Two-way aggregate fact table: $5 \times 40 \times 30 \times 3 \times 12$
 $: 216000$ rows

One-way aggregate fact table: $3 \times 365 \times 10,000 \times 40 \times 30$
 $: 6570000$ rows

Star Schema

→

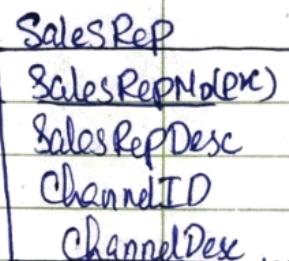
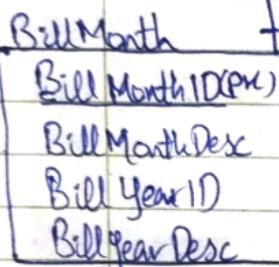
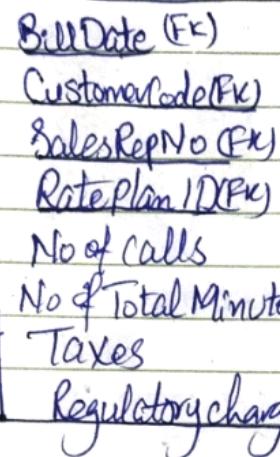
Bill Date



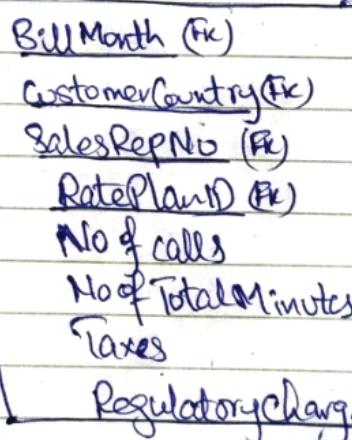
BillingBaseFactTable

(10)

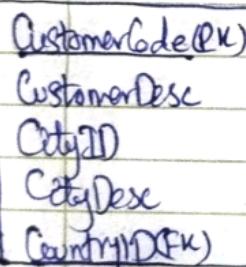
Rate Plan



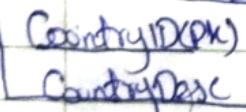
Two-way aggregate fact table



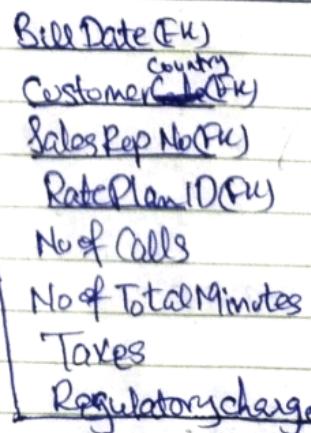
Customer



Customer Country



One-way aggregate fact table



Q3

(1c)

- ① Dependent data marts directly create data from centralized data warehouse. They ensure data consistency across enterprise. Whereas, independent data marts are standalone systems that create data from various resources without relying on a central data warehouse and hence result in a data inconsistency issues. There are various architectures like hub and spoke, federated, datamarts, Centralized datawarehouse and independent data mart.

- ② Hence result in a data inconsistency issues. There are various architectures like hub and spoke, federated, datamarts, Centralized datawarehouse and independent data mart.

- ③ Snowflake schema is normalized version of star schema where dimension tables are split into multiple related tables whereas, star schema is a simpler structure where central base fact table is connected to denormalized dimension tables.

Advantage: ① Minimal data redundancy - ② Storage is reduced.

- ④ Pre-join denormalization is a technique involving joining tables in advance and storing results in a single table to improve query performance. Whereas, column replication is a technique of replicating columns across multiple tables to avoid joins and simplified queries for better query performance. Since you do not need to join tables at query time, so it improves overall performance.

- ⑤ Partitioned cubes are large OLAP cubes that are divided into smaller, manageable portions for performance optimization. They are useful for larger data sets allowing faster query processing and improved load performance.

- ⑥ There are 4 different types of OLAP:

MOLAP

ROLAP

DOLAP

HOLAP

Since ROLAP stores data in relational databases and have slower query performance as it generates query dynamically but has better scalability on large datasets.

Data Warehousing and Business Intelligence (DS3003)

Date: November 4th 2024

Course Instructor

M. Ishaq Raza

Sessional-II Exam

Total Time (Hrs.): 1
Total Marks: 25
Total Questions: 3

Roll No

Section

Student Signature

Do not write below this line.

Note: Please ensure that you attempt all questions and their respective parts in the given order.
You may use a calculator.

CLO # 1: Demonstrate an appreciation of the role that DW and BI play in enhancing the decision-making process.

Q. No 1: Give the appropriate answers to the following questions: [3+2]

- When it is appropriate to use full data refresh loading strategy? Also suggest at least two practical steps that expedite the full data refresh loading process.
- Briefly explain the two most efficient immediate data extraction techniques.

Consider the following description for the next Questions:

Consider the following tables and statistics which are part of a Library system:

Book (BookID, Title, Publisher, PublishYear, Author, ...);

BookLoan (BookID, BranchID, CardNo, DateOut, DueDate, ...);

Assume Book and BookLoan tables containing 100,000 and 2,000,000 rows respectively. Each table row and each index entry take 200 bytes and 15 bytes of space respectively. Data block size is 8KB and available memory size is 100 blocks. Suppose selectivity of publisher 'Wiley' = 10%, publishYear '2023' = 8%, author 'Inmon' = 4%, and author 'Kimball' = 2%.

CLO # 1: Demonstrate an appreciation of the role that DW and BI play in enhancing the decision-making process.

Q. No 2: Calculate the total I/O cost for the Query using the following indexed access paths. Show all steps clearly. [10]

Query: `SELECT * FROM book
WHERE publisher= 'Wiley' AND publishYear= '2023' AND author IN ('Inmon', 'Kimball');`

- Clustered Index access (Assume clustered index exit on Author column)
- Static Bitmap Index access (Assume static bitmap indexes exist on Publisher, PublishYear, and Author columns separately)

National University of Computer and Emerging Sciences 0027
Lahore Campus

CLO # 1: Demonstrate an appreciation of the role that DW and BI play in enhancing the decision-making process.

Q. No 3: Calculate the total I/O cost to execute this Query using the following joining techniques. Show all steps clearly. Assume static bitmap indexes exist on Publisher, PublishYear, and Author columns of Book table and there is no index on BookLoan table. [3+3+4= 10]

Query: *SELECT * FROM book JOIN BookLoan ON Book.BookID = BookLoan.BookID
WHERE publisher= 'Wiley' AND publishYear= '2023' AND author IN ('Inmon', 'Kimball');*

- a. Sort Merge Join
- b. Hash Join
- c. Nested Loop Join (*Identify the most efficient variant of NLJ in this scenario, then compute the I/O cost of that variant only.*)

Student's Name Ameer Abdullah Signature Qut

Roll No. 21L-S694 Section BDS-STA Date 4th Nov, 2024

Q1)

a) we use full state refresh loading strategy when our large data changes occur - like the ratio between existing and new data size exceeds 10%. The steps to complete the full state refresh is
 i) Batch processing
 ii) Parallel processing
 iii) Drop and rebuild indexes.

b) i) Transaction log: They are already in relational database management and does not require extensive customization. It is one of the most cost-effective technique.

ii) Database Triggers: They are built in database and automatically track changes without requiring any additional tools or processes. They are easy to setup and manage.

Q2)	$r_B = 100000$	$bfr = \frac{8 \times 1024}{200} = 40$ blocks
	$r_{B2} = 2000000$	
	$R = 200$ bytes	$bfr_i = \frac{8 \times 1024}{15} = 546$ index entries
	$R_i = 15$ bytes	
	$B = 8KB \rightarrow 8 \times 1024$	$b_B = \frac{100000}{40} = 2500$ blocks
	$K = 100$ blocks	
Publisher	$Wiley = 10\%$	
PublishYear	$2023 = 8\%$	$b_{BL} = \frac{2000000}{40} = 50000$ blocks
Author	$Tanen = 4\%$	
Author	$Kimball = 2\%$	
		$1 : 20$

Q / Part No.

$$b_{BI} = \left\lceil \frac{100000}{546} \right\rceil = 184 \text{ index blocks}$$

$$b_{BLI} = \left\lceil \frac{2000000}{546} \right\rceil = 3664 \text{ index blocks}$$

(10)

$$\text{Qualifying rows} = 10\% \times 8\% \times 6\% \times 100000 \Rightarrow 48$$

$$\text{Qualifying blocks: } \left\lceil \frac{48}{40} \right\rceil = 2 \text{ blocks}$$

(Q2)

(a) Clustered Index access :

$$6\% \times 100000 = 6000$$

$$\text{Base table access cost : } \left\lceil \frac{6000}{40} \right\rceil = 150$$

$$\text{Index table access cost : } \left\lceil \frac{6000}{546} \right\rceil = 11$$

$$\text{Total Cost (Dense Index)} : 150 + 11 = 161 \text{ I/O}$$

$$\rightarrow \text{Sparse Index} \rightarrow \text{Index table access cost : } 1 + 1 = 2$$

$$\rightarrow 150 + 2 = 152 \text{ I/O}$$

$$(b) \left\lceil \frac{100000}{8 \times 15 \times 546} \right\rceil = 2 \text{ (Size for one bitmap)}$$

$$(2 + 2 + 2 + 2) = 8$$

$$\begin{aligned} \text{Since } QR &\perp b_B \\ 48 &\perp 2800 \end{aligned}$$

5

$$\text{So, Total Cost : } 8 + 48 = 56 \text{ I/O}$$

Q3

(c)

QB \leftarrow K] so we do not merge sort Book table
 SMJ: 2 \leftarrow 100]

$FTS + (BookSort) + (BookLoanSort) + (MergeCost)$

$$2500 + 2 + \left(50000 \times \log_2 \left(\frac{50000}{100} \right) \right) + (2 + 50000)$$

$$2500 + 2 + 282193 + 50002$$

$$\text{Total I/O cost : } 334697 \text{ I/O}$$

we took $FTS = 56$ due to static bitmap index on outer table -

(b) Hash Join: NO Partitioning cost for both table as smaller table can fit into the memory -

$$FTS + Hashing = \text{Total I/O cost}$$

$$2500 + (2 + 50000) = 52502 \text{ I/O}$$

(c) Nested Loop Join:

Best Variant:

Clustered NLJ: $FTS + QR \times (\text{Index access cost} + \text{base table access cost})$

$$: 2500 + 48 \times (1 + 1) = 2596 \text{ I/O}$$

~~no index exist on inner table!~~

2nd Best Variant:

$$\text{Indexed NLJ: } 2500 + 48 \times (1 + 20) = 3508 \text{ I/O}$$

→ Since Basic And Block NLJ will be highly costly so those are not good option -

→ ALSO as we can see that Indexed NLJ cost is significantly greater so we consider clustered NLJ as the best -

~~$* \text{ AS bitmap indexes exist } 56 + 48(1+1) = 152 \text{ I/O } *$~~

(c) $56 + (2 \times 50000)$ correct ✓

Data Warehousing and Business Intelligence (DS3003)

Date: December 28th 2024

Course Instructor

M. Ishaq Raza

Final Exam

Total Time (Hrs.): 3
Total Marks: 60
Total Questions: 5

Roll No

Section

Student Signature

Do not write below this line.

Note: Please ensure that you attempt all questions and their respective parts in the given order.
You may use a calculator.

CLO # 2: Demonstrate an understanding of the fundamental concepts of the Star and the Snowflake Schema; learn how to design the schema of a DW based on these two models.

Q. No 1: [12+4= 16]

- As the data design specialist for the data warehouse project team of an eCommerce system, **design a star schema** that includes a base fact table with at least four dimensions and three aggregate fact tables (a 3-way, 2-way, and 1-way aggregation). Provide the possible attributes for each dimension and fact table. Show the primary keys, foreign keys and all the relationships between the dimension and fact tables. Additionally, design a dimension table that supports the preservation of historical changes. Note: Draw a single diagram that includes the base fact table and the aggregate fact tables.
- Take appropriate cardinality (i.e., number of rows) of each of the above dimension and their levels. Estimate the potential size (in number of rows) of the above base fact table as well as aggregate fact tables.

Product ID Total Sales
Sale ID Customer ID Taxes
Customer ID Billing Date Products sold/overtime
Products sold/overtime No. of customers

CLO # 1: Demonstrate an appreciation of the role that DW and BI play in enhancing the decision-making process.

Q. No 2: Give the appropriate answers to the following questions: [15]

- Briefly discuss the three major types of architectures for building a data warehouse.
- Give at least three reasons why you think ETL functions are most challenging in a data warehouse environment.
- Drill-down, drill-across, and drill-through are analytical operations that can be performed on an OLAP cube. Provide a brief description of each.
- Discuss at least three advantages of using materialized views.
- Can we use a combination of vertical and horizontal partition techniques together for optimal performance? Illustrate your answer with an example.

$$\frac{2\% \times 8\% \times 20000}{819} = 1+320$$

$$\frac{2\% \times 8\% \times 20000}{819} = 1+240$$

$$\frac{4\% \times 8\% \times 20000}{819} = 1+640$$

$$\frac{4\% \times 6\% \times 20000}{819} = 1+480$$

$$819 \rightarrow 1684$$

CLO # 1: Demonstrate an appreciation of the role that DW and BI play in enhancing the decision-making process.

Q. No 3: Suppose you have the following market basket data. [5]

TID	Items-Bought
1	{A, C, F}
2	{B, C, D}
3	{A, B, C, D}
4	{B, D}
5	{E}

Find all frequent itemsets using Apriori algorithm with min_sup=2, i.e., any itemset occurring in less than 2 transactions is infrequent. Also list all the strong association rules with min_sup=2 and min_conf=100%.

$$\begin{array}{l} BC \rightarrow D \\ CD \rightarrow B \\ CD \rightarrow C \\ A \rightarrow C \\ A \rightarrow B \\ B \rightarrow D \\ B \rightarrow C \\ D \rightarrow B \end{array}$$

$$\begin{array}{l} BC \rightarrow D \\ CD \rightarrow B \\ CD \rightarrow C \\ A \rightarrow C \\ B \rightarrow D \\ D \rightarrow B \end{array}$$

Consider the following description for the next two Questions:

Consider the following tables and statistics which are part of a vehicle sales system:

Vehicle (VehicleID, Make, Model, Color, ...);

Sales (SalesID, VehicleID, SalesPersonID, SalesDate, Price, ...);

Assume Vehicle and Sales tables containing 200,000 and 150,000 rows respectively. Each table row and each index entry take 100 bytes and 20 bytes of space respectively. Data block size is 16KB and available memory size is 50 blocks. Suppose selectivity of Model '2024'= 2%, Model '2023'= 4%, Color 'White'=8, and Color 'Black'=6%.

CLO # 1: Demonstrate an appreciation of the role that DW and BI play in enhancing the decision-making process.

Q. No 4: Calculate the total I/O cost to execute this Query using the following indexed access paths. Show all steps clearly. [12]

Query: `SELECT * FROM vehicle
WHERE model IN (2024, 2023) AND color IN ('White', 'Black');`

- Single Index access (Assume single indexes exist on model and color columns separately) 1242
- Dynamic Bitmap Index access (Assume single indexes exist on model and color columns separately) 1277
- Composite Index Access (Assume a composite index exist on model and color columns, with index entry size=20 bytes) 1684

CLO # 1: Demonstrate an appreciation of the role that DW and BI play in enhancing the decision-making process.

Q. No 5: Calculate the total I/O cost to execute this Query using the following joining techniques. Show all steps clearly. Assume there is no index exist on any table. [12]

Query: `SELECT * FROM vehicle JOIN sales ON vehicle.vehicleID = sales.vehicleID
WHERE model IN (2024, 2023) AND color IN ('White', 'Black');`

- Hash Join 2159 or 2156
- Sort Merge Join 6042
- Nested Loop Join (Identify the most efficient variant of NLJ in this scenario, then compute the I/O cost of that variant only.) 11358