

ASSESSMENT SUBMISSION	
Module Title:	Advanced Databases
Module Code:	KL7011
Academic Year / Semester:	2023-24 / Semester 1
Module Tutor / Email (all queries):	Akhtar Ali akhtar.ali@northumbria.ac.uk
% Weighting (to overall module):	40%
Assessment Title:	Assignment 2: team-work
Group Work	This assessment is designed to be undertaken by a group comprising TWO students. If you cannot find someone to work with then you can do the assessment all by yourself.
Date of Handout to Students:	28 th November 2023
Mechanism for Handout:	Module Blackboard Site & Live Session in Week 9
Deadline for Submission Attempt by Students:	22 nd January 2024 @ 23:59 GMT
Mechanism for Submission:	Document upload to Module Blackboard Site
Submission Format / Word Count	Please upload your written report as a single PDF document
Date by which Work, Feedback and Marks will be returned:	26 th February 2024
Mechanism for return of Feedback and Marks:	Mark and written feedback will be uploaded to the Module Site on Blackboard. For further queries please email module tutor.
Student IDs/Oracle Username (DWUs and DMU)	W22058786, W23008776 / DWU35, DMU17
Names of students in the Group	Christy Joseph, Muhammad Shibil Alambilatt Keloth
Group No	17

Instructions on Assessment:

- ♦ You are expected to produce a word-processed answer to this assignment. Please use Arial font and a font size of 12 for text. For SQL code and output, you can use courier new font and a minimum size of 10, which preserves SQL format and layout. Where necessary, screenshots of SQL output may be used instead of plain text.
- ♦ You are required to use the Harvard Style of referencing and citation. The “*Cite them right*” guide is recommended for *referencing and citation* (Pears and Shields, 2008) which should be followed throughout your answer especially Part 3. Please do not include references to lecture notes.
- ♦ **ONLY ONE submission is required for each group to be submitted on Blackboard.**
- ♦ The names of students in the group must be provided and must match with the group no and names already agreed on the shared document.
- ♦ Marks allocated for your submission will be shared equally by all the students within the group (a max of 2 members per group). However, if some members have not contributed to the assignment as agreed and expected of them, then a peer-assessment form should be filled and submitted on the Blackboard by each member of the group. See Appendixes 3, 4 and 5.

Personalising your SQL output/prompt

Before executing any **SQL code** for this assignment, you should personalise your SQL output / prompt by running SET SQLPROMPT “DWUn >”, i.e., *double-quote* followed by your Data Warehouse user name (which could be one of the two members of the group) followed by > and then a *space* and *double-quote* as shown in the screenshot below. Likewise, for Part 2, you must personalise the SQL prompt using your DMU username linked to your group.

```
Select SQL Plus

SQL> SET SQLPROMPT "DWU152 > "
DWU152 > desc sales
Name                                     Null?    Type
-----
QUANTITY_SOLD                          NOT NULL NUMBER(3)
AMOUNT_SOLD                            NOT NULL NUMBER(10,2)
PROD_ID                                NOT NULL NUMBER(6)
CUST_ID                                NOT NULL NUMBER
TIME_ID                                NOT NULL DATE
CHANNEL_ID                              NOT NULL CHAR(1)
PROMO_ID                                NOT NULL NUMBER(6)

DWU152 > desc channels
Name                                     Null?    Type
-----
CHANNEL_ID                              NOT NULL CHAR(1)
CHANNEL_DESC                            NOT NULL VARCHAR2(20)
CHANNEL_CLASS                           NOT NULL VARCHAR2(20)

DWU152 >
```

Assignment Questions

Part 1: Data Warehousing Tasks (50 Marks)

This part is based on the **Sales History** scenario as described in Appendix 1.

You must submit all the SQL queries and any other code that you wrote in answering any of the tasks / questions (e.g., the use of EXPLAIN PLAN statements for the queries and their outputs using Spooling or other suitable means).

- (A) Study the index definitions in `sh_idx.sql`. These indexes have already been created in SH2. Whatever indexes you decide to create for this task should be the result of your own research and thinking, and be different than those already exist in SH2 or those indexes defined in the Oracle Data Warehousing Guide (Potineni, 2021) or those of other students.

You need to design *two* queries such that each query involves at least *three* different tables and at least *one* aggregate function. You need to ensure that your queries have adequate *selectivity* such that if suitable indexes were available in your DWU version of the database, the queries would have performed more efficiently.

You need to identify and justify at least two indexes to improve the performance of your queries. Then create your proposed indexes in your DWU version of the database. You need to run your queries before and after creating your proposed indexes and report EXPLAIN PLAN outputs and make sure that your proposed indexes have been used by your queries and have improved their performance significantly.

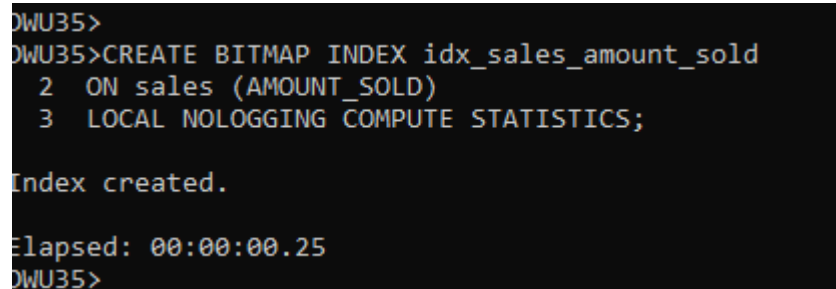
Then critically discuss the differences in the performance of your queries with and without the proposed indexes. You need to critically review and cite relevant database literature to support your choice of indexes and how you dealt with the issue of selectivity in your queries.

(20 marks)

Answer Part 1 (A)

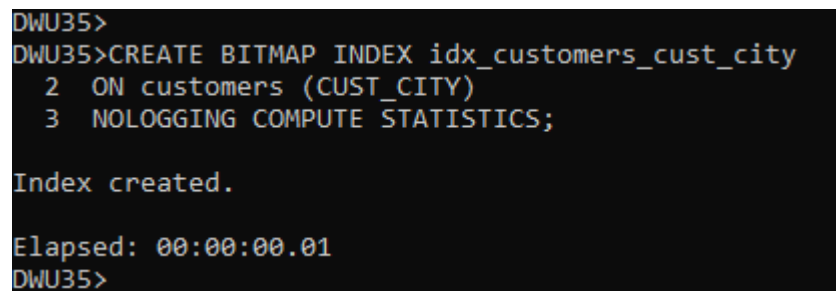
Provide the SQL code and output of the 2 indexes you have created on your DWU database for comparing their performance impact on (i.e., these indexes must not exist in SH2) (4 Mark). **Make sure the SQL code you provide is plain text and the output is a screenshot.**

```
CREATE BITMAP INDEX idx_sales_amount_sold  
ON sales (AMOUNT_SOLD)  
LOCAL NOLOGGING COMPUTE STATISTICS;
```



```
DWU35>  
DWU35>CREATE BITMAP INDEX idx_sales_amount_sold  
2 ON sales (AMOUNT_SOLD)  
3 LOCAL NOLOGGING COMPUTE STATISTICS;  
  
Index created.  
  
Elapsed: 00:00:00.25  
DWU35>
```

```
CREATE BITMAP INDEX idx_customers_cust_city  
ON customers (CUST_CITY)  
NOLOGGING COMPUTE STATISTICS;
```



```
DWU35>  
DWU35>CREATE BITMAP INDEX idx_customers_cust_city  
2 ON customers (CUST_CITY)  
3 NOLOGGING COMPUTE STATISTICS;  
  
Index created.  
  
Elapsed: 00:00:00.01  
DWU35>
```

Provide the rationale and justification of creating the above indexes based on your own research and citing appropriate literature here and providing references in the "References and Bibliography" section at the end of the report (4 Marks):

In the provided scenario, two bitmap indexes have been created on the "sales" and "customers" tables. Let's analyze the rationale and justification for each index based on general principles and literature:

1. Bitmap Index on "sales" Table (idx_sales_amount_sold):

- Column Indexed: AMOUNT_SOLD

- Rationale:

- The "sales" table likely involves queries that filter and aggregate data based on the sales amount.

- Bitmap indexes are particularly effective for columns with low cardinality (a small number of distinct values), such as status flags or categories. If the "AMOUNT_SOLD" column has a limited range of values, a bitmap index can efficiently speed up queries that involve filtering by sales amount. (O'Neil,1997)

- Justification:

- Bitmap indexes are well-suited for scenarios where queries involve range or equality conditions on low cardinality columns. The index provides a compact representation of the data distribution, enabling faster retrieval.

- The NOLOGGING option reduces the amount of redo log generated during index creation, which can be beneficial for performance in certain situations, such as bulk data loading.(Oracle)

2. Bitmap Index on "customers" Table (idx_customers_cust_city):

- Column Indexed: CUST_CITY
- Rationale:
 - The "customers" table likely involves queries that filter and group data based on the customer's city.
 - Bitmap indexes are suitable for columns with low cardinality, and if the "CUST_CITY" column has a limited set of distinct values, a bitmap index can significantly improve query performance.(Oracle)
- Justification:
 - Bitmap indexes are efficient for handling queries with equality conditions on low cardinality columns, such as searching for customers in a specific city.
 - The NOLOGGING option may be chosen to minimize the impact on transaction logging during the index creation process, which can be crucial for performance in data-intensive environments.(Oracle)

Provide the 2 SQL queries you are going to run to compare the performance impact of the above 2 Indexes on DWU (4 marks). Make sure the SQL code you provide is plain text.

```
SELECT
    C.CUST_CITY,
    P.PROD_CATEGORY,
    SUM(S.AMOUNT_SOLD) AS TOTAL_AMOUNT_SOLD
FROM
    CUSTOMERS C
    JOIN SALES S ON C.CUST_ID = S.CUST_ID
    JOIN PRODUCTS P ON S.PROD_ID = P.PROD_ID
WHERE
    C.CUST_CITY IS NOT NULL
```

```
        AND C.CUST_CITY IN (
            SELECT CUST_CITY
            FROM CUSTOMERS
            WHERE CUST_CITY IS NOT NULL
            GROUP BY CUST_CITY
            HAVING COUNT(*) > 5
        )
        AND S.AMOUNT_SOLD < 50
GROUP BY
    C.CUST_CITY, P.PROD_CATEGORY;
```

```
DWU35>SELECT
 2      C.CUST_CITY,
 3      P.PROD_CATEGORY,
 4      SUM(S.AMOUNT_SOLD) AS TOTAL_AMOUNT_SOLD
 5  FROM
 6      CUSTOMERS C
 7      JOIN SALES S ON C.CUST_ID = S.CUST_ID
 8      JOIN PRODUCTS P ON S.PROD_ID = P.PROD_ID
 9  WHERE
10      C.CUST_CITY IS NOT NULL
11      AND C.CUST_CITY IN (
12          SELECT CUST_CITY
13          FROM CUSTOMERS
14          WHERE CUST_CITY IS NOT NULL
15          GROUP BY CUST_CITY
16          HAVING COUNT(*) > 5
17      )
18      AND S.AMOUNT_SOLD < 50
19  GROUP BY
20      C.CUST_CITY, P.PROD_CATEGORY;
```

```
SELECT
    C.CUST_CITY,
    C.CUST_GENDER,
    P.PROD_CATEGORY,
    COUNT(*) AS CUSTOMER_COUNT,
    SUM(S.AMOUNT_SOLD) AS TOTAL_AMOUNT_SOLD
FROM
    CUSTOMERS C
    JOIN SALES S ON C.CUST_ID = S.CUST_ID
    JOIN PRODUCTS P ON S.PROD_ID = P.PROD_ID
WHERE
    C.CUST_CITY IS NOT NULL
    AND C.CUST_CITY LIKE 'A%' OR C.CUST_CITY LIKE 'Sa%'
    AND S.AMOUNT_SOLD < 50
GROUP BY
    C.CUST_CITY, C.CUST_GENDER, P.PROD_CATEGORY;
```



```
DWU35>
DWU35>SELECT
  2     C.CUST_CITY,
  3     C.CUST_GENDER,
  4     P.PROD_CATEGORY,
  5     COUNT(*) AS CUSTOMER_COUNT,
  6     SUM(S.AMOUNT_SOLD) AS TOTAL_AMOUNT_SOLD
  7 FROM
  8     CUSTOMERS C
  9     JOIN SALES S ON C.CUST_ID = S.CUST_ID
 10     JOIN PRODUCTS P ON S.PROD_ID = P.PROD_ID
 11 WHERE
 12     C.CUST_CITY IS NOT NULL
 13     AND C.CUST_CITY LIKE 'A%' OR C.CUST_CITY LIKE 'Sa%'
 14     AND S.AMOUNT_SOLD < 50
 15 GROUP BY
 16     C.CUST_CITY, C.CUST_GENDER, P.PROD_CATEGORY;_
```

Provide Explain Plan statements & outputs for the above 2 SQL queries you have run to compare the performance impact of your 2 indexes on DWU before and after creating your proposed indexes (4 marks). Make sure the SQL code you provide is plain text and the output is a screenshot.

CUST_CITY	PROD_CATEGORY	TOTAL_AMOUNT_SOLD
Courchevel	Men	120.2
Bedford	Women	269
Heidelberg	Women	650.2
Potsdam	Women	210.1
Offenbach	Men	193
Thame	Men	97
Saint-Brieuc	Girls	526.8
Veldhoven	Men	215.6
Cochin Kochi	Women	185.2
Mandelieu la Napoule	Boys	70
Alleppey	Men	135

CUST_CITY	PROD_CATEGORY	TOTAL_AMOUNT_SOLD
Woodstock	Boys	255
Mc Kean	Women	208
Warsaw	Girls	32
Bristol	Women	60
Sainte-Croix-du-Mont	Girls	48
Zandvoort	Girls	162
Relecq-Kerhuon	Men	91
Toulouse	Men	175
Groesbeek	Men	72
Woodstock	Men	38

1748 rows selected.
Elapsed: 00:00:01.31
DWU35>

The time taken for execution before indexing was 01.31 seconds for this first query.

```
EXPLAIN PLAN FOR
SELECT
  C.CUST_CITY,
  P.PROD_CATEGORY,
  SUM(S.AMOUNT_SOLD) AS TOTAL_AMOUNT_SOLD
FROM
```

```
CUSTOMERS C
JOIN SALES S ON C.CUST_ID = S.CUST_ID
JOIN PRODUCTS P ON S.PROD_ID = P.PROD_ID
WHERE
  C.CUST_CITY IS NOT NULL
  AND C.CUST_CITY IN (
    SELECT CUST_CITY
    FROM CUSTOMERS
    WHERE CUST_CITY IS NOT NULL
    GROUP BY CUST_CITY
    HAVING COUNT(*) > 5
  )
  AND S.AMOUNT_SOLD < 50
GROUP BY
  C.CUST_CITY, P.PROD_CATEGORY;
```

```
SELECT * FROM TABLE(DBMS_XPLAN.DISPLAY);
```

```
DWU35>
DWU35>EXPLAIN PLAN FOR
 2  SELECT
 3    C.CUST_CITY,
 4    P.PROD_CATEGORY,
 5    SUM(S.AMOUNT_SOLD) AS TOTAL_AMOUNT_SOLD
 6  FROM
 7    CUSTOMERS C
 8    JOIN SALES S ON C.CUST_ID = S.CUST_ID
 9    JOIN PRODUCTS P ON S.PROD_ID = P.PROD_ID
10 WHERE
11   C.CUST_CITY IS NOT NULL
12   AND C.CUST_CITY IN (
13     SELECT CUST_CITY
14     FROM CUSTOMERS
15     WHERE CUST_CITY IS NOT NULL
16     GROUP BY CUST_CITY
17     HAVING COUNT(*) > 5
18   )
19   AND S.AMOUNT_SOLD < 50
20 GROUP BY
21   C.CUST_CITY, P.PROD_CATEGORY;
```

Explained.

Elapsed: 00:00:00.01

DWU35>

DWU35>

```
DWU35>SELECT * FROM TABLE(DBMS_XPLAN.DISPLAY);
```

```

DWU35>SELECT * FROM TABLE(DBMS_XPLAN.DISPLAY);

PLAN_TABLE_OUTPUT
-----
Plan hash value: 3870217569

-----
| Id | Operation                      | Name      | Rows  | Bytes | Cost (%CPU)| Time     | Pstart | Pstop |
-----+-----+-----+-----+-----+-----+-----+-----+-----+
|  0 | SELECT STATEMENT                |           |      8 | 4400  | 3950  (1) | 00:00:01 |        |       |
|  1 |   HASH GROUP BY                  |           |      8 | 4400  | 3950  (1) | 00:00:01 |        |       |
| * 2 |    HASH JOIN                     |           |     162 | 8100  | 3949  (1) | 00:00:01 |        |       |
| * 3 |      HASH JOIN                   |           |     162 | 6318  | 3846  (1) | 00:00:01 |        |       |
|  4 |        PARTITION RANGE ALL       |           |    3247 | 45458 | 3295  (1) | 00:00:01 |      1 |     17 |
| * 5 |          TABLE ACCESS FULL      | SALES     |    3247 | 45458 | 3295  (1) | 00:00:01 |      1 |     17 |
-----
PLAN_TABLE_OUTPUT
-----
| * 6 |    HASH JOIN RIGHT SEMI          |           |    2500 | 62500 | 552   (2) | 00:00:01 |        |       |
|  7 |      VIEW                        | VW_NSO_1  |        31 | 310   | 277   (2) | 00:00:01 |        |       |
| * 8 |        HASH GROUP BY             |           |        31 | 310   | 277   (2) | 00:00:01 |        |       |
|  9 |          TABLE ACCESS FULL      | CUSTOMERS |    50000 | 488K   | 274   (1) | 00:00:01 |        |       |
| 10 |            TABLE ACCESS FULL    | CUSTOMERS |    50000 | 732K   | 274   (1) | 00:00:01 |        |       |
| 11 |              TABLE ACCESS FULL  | PRODUCTS |    10000 | 107K   | 102   (0) | 00:00:01 |        |       |
-----

Predicate Information (identified by operation id):
-----
PLAN_TABLE_OUTPUT
-----
      2 - access("S"."PROD_ID"="P"."PROD_ID")
      3 - access("C"."CUST_ID"="S"."CUST_ID")
      5 - filter("S"."AMOUNT_SOLD"<50)
      6 - access("C"."CUST_CITY"="CUST_CITY")
      8 - filter(COUNT(*)>5)

27 rows selected.

Elapsed: 00:00:00.03
DWU35>_

```

```

DWU35>SELECT * FROM TABLE(DBMS_XPLAN.DISPLAY);

PLAN_TABLE_OUTPUT
-----
Plan hash value: 2837496622

-----
| Id | Operation                                | Name                                | Rows  | Bytes | Cost (%CPU)| Time     | Pstart | Pstop |
-----+-----+-----+-----+-----+-----+-----+-----+-----+
|  0 | SELECT STATEMENT                        |                                     |      8 | 4400 | 1039 (1)| 00:00:01 |        |       |
|  1 | HASH GROUP BY                          |                                     |      8 | 4400 | 1039 (1)| 00:00:01 |        |       |
|* 2 | HASH JOIN                              |                                     |    162 | 8100 | 1038 (1)| 00:00:01 |        |       |
|* 3 | HASH JOIN                              |                                     |    162 | 6318 | 935 (1)| 00:00:01 |        |       |
|  4 | PARTITION RANGE ALL                    |                                     |    3247 | 45458 | 643 (0)| 00:00:01 |        |       |
|  5 | TABLE ACCESS BY LOCAL INDEX ROWID BATCHED | SALES                             |    3247 | 45458 | 643 (0)| 00:00:01 |        |       |
-----+-----+-----+-----+-----+-----+-----+-----+
PLAN_TABLE_OUTPUT
-----
|  6 | BITMAP CONVERSION TO ROWIDS            |                                     |      8 | 4400 | 1039 (1)| 00:00:01 |        |       |
|* 7 | BITMAP INDEX RANGE SCAN                | IDX_SALES_AMOUNT_SOLD             |    2500 | 62500 | 292 (1)| 00:00:01 |        |       |
|* 8 | HASH JOIN RIGHT SEMI                   |                                     |    31 | 310 | 18 (0)| 00:00:01 |        |       |
|  9 | VIEW                                    | VW_NS0_1                          |    31 | 310 | 18 (0)| 00:00:01 |        |       |
|* 10 | HASH GROUP BY                         |                                     |    31 | 310 | 18 (0)| 00:00:01 |        |       |
| 11 | BITMAP CONVERSION TO ROWIDS            |                                     |   50000 | 488K | 18 (0)| 00:00:01 |        |       |
| 12 | BITMAP INDEX FULL SCAN                 | IDX_CUSTOMERS_CUST_CITY            |   50000 | 732K | 274 (1)| 00:00:01 |        |       |
| 13 | TABLE ACCESS FULL                     | CUSTOMERS                          |   50000 | 732K | 274 (1)| 00:00:01 |        |       |
| 14 | TABLE ACCESS FULL                     | PRODUCTS                           |  100000 | 107K | 102 (0)| 00:00:01 |        |       |
-----+-----+-----+-----+-----+-----+-----+
PLAN_TABLE_OUTPUT
-----
Predicate Information (identified by operation id):
-----
   2 - access("S"."PROD_ID"="P"."PROD_ID")
   3 - access("C"."CUST_ID"="S"."CUST_ID")
   7 - access("S"."AMOUNT_SOLD"<=0)
     filter("S"."AMOUNT_SOLD"<=0)
   8 - access("C"."CUST_CITY"="CUST_CITY")
  10 - filter(COUNT(*)>5)

31 rows selected.

Elapsed: 00:00:00.07
DWU35>

```

After the indexes were created the explain plan table shows that the indexes were used in calling the data for the first query.

```

CUST_CITY          PROD_CATEGORY          TOTAL_AMOUNT_SOLD
-----
Courchevel         Men                      120.2
Bedford            Women                    269
Heidelberg         Women                    650.2
Potsdam            Women                    210.1
Offenbach          Men                      193
Thame              Men                      97
Saint-Brieuc       Girls                    526.8
Veldhoven          Men                      215.6
Cochin Kochi       Women                    185.2
Mandelieu la Napoule Boys                      70
Alleppey           Men                      135

CUST_CITY          PROD_CATEGORY          TOTAL_AMOUNT_SOLD
-----
Woodstock          Boys                     255
Mc Kean            Women                    208
Warsaw             Girls                     32
Bristol            Women                     60
Sainte-Croix-du-Mont Girls                     48
Zandvoort          Girls                    162
Relecq-Kerhuon     Men                       91
Toulouse           Men                     175
Groesbeek          Men                       72
Woodstock          Men                       38

1748 rows selected.

Elapsed: 00:00:00.49
DWU35>

```

And the time taken for the execution of the code is down to 0.49 second.

CUST_CITY	C PROD_CATEGORY	CUSTOMER_COUNT	TOTAL_AMOUNT_SOLD
Adelaide	F Boys	241	72149.95
Alma	M Girls	173	70706.8
Almere	M Girls	288	101547.5
Aix-les-Bains	M Men	232	254912.4
Santos	M Men	8	329.3
San Mateo	F Women	5	222
San Nicolas	M Men	11	437.3
Alleppey	F Boys	186	59432
Amsterdam	F Boys	37	6635.8
Alkmaar	F Boys	6	1705
Albion	M Men	35	19050
Aix-en-Provence	M Men	3	9104
San Nicolas	F Girls	2	92

288 rows selected.
Elapsed: 00:00:00.37
DWU35>

The time taken for execution of the second query before indexing was 0.37 seconds.

```
EXPLAIN PLAN FOR
SELECT
    C.CUST_CITY,
    C.CUST_GENDER,
    P.PROD_CATEGORY,
    COUNT(*) AS CUSTOMER_COUNT,
    SUM(S.AMOUNT_SOLD) AS TOTAL_AMOUNT_SOLD
FROM
    CUSTOMERS C
    JOIN SALES S ON C.CUST_ID = S.CUST_ID
    JOIN PRODUCTS P ON S.PROD_ID = P.PROD_ID
WHERE
    C.CUST_CITY IS NOT NULL
    AND C.CUST_CITY LIKE 'A%' OR C.CUST_CITY LIKE 'Sa%'
    AND S.AMOUNT_SOLD < 50
GROUP BY
    C.CUST_CITY, C.CUST_GENDER, P.PROD_CATEGORY;

SELECT * FROM TABLE(DBMS_XPLAN.DISPLAY);
```

```
DWU35>SELECT * FROM TABLE(DBMS_XPLAN.DISPLAY);
```

PLAN_TABLE_OUTPUT

Plan hash value: 2697944916

Id	Operation	Name	Rows	Bytes	TempSpc	Cost (%CPU)	Time	Pstart	Pstop
0	SELECT STATEMENT		2480	101K		5011 (1)	00:00:01		
1	HASH GROUP BY		2480	101K		5011 (1)	00:00:01		
* 2	HASH JOIN		70316	2884K		5006 (1)	00:00:01		
3	TABLE ACCESS FULL	PRODUCTS	10000	107K		102 (0)	00:00:01		
* 4	HASH JOIN		70316	2128K	1424K	4903 (1)	00:00:01		
5	TABLE ACCESS FULL	CUSTOMERS	50000	830K		274 (1)	00:00:01		

PLAN_TABLE_OUTPUT

6	PARTITION RANGE ALL		1016K	13M		3294 (1)	00:00:01	1	17
7	TABLE ACCESS FULL	SALES	1016K	13M		3294 (1)	00:00:01	1	17

Predicate Information (identified by operation id):

```

2 - access("S"."PROD_ID"="P"."PROD_ID")
4 - access("C"."CUST_ID"="S"."CUST_ID")
    filter("C"."CUST_CITY" LIKE 'A%' OR "C"."CUST_CITY" LIKE 'Sa%' AND "S"."AMOUNT_SOLD"<50)

```

21 rows selected.

Elapsed: 00:00:00.02
DWU35>

Output after indexing:

```

DWU35>SELECT * FROM TABLE(DBMS_XPLAN.DISPLAY);

PLAN_TABLE_OUTPUT
-----
Plan hash value: 1114399456

-----
| Id | Operation | Name | Rows | Bytes | Cost (%CPU) | Time | Pstart | Pstop |
-----
| 0 | SELECT STATEMENT | | 6 | 288 | 4594 (3) | 00:00:01 | | |
| 1 | HASH GROUP BY | | 6 | 288 | 4594 (3) | 00:00:01 | | |
| * 2 | HASH JOIN | | 1019K | 46M | 4523 (1) | 00:00:01 | | |
| 3 | TABLE ACCESS FULL | PRODUCTS | 10000 | 107K | 102 (0) | 00:00:01 | | |
| 4 | VIEW | VW_3F_SET$2263B486 | 1019K | 35M | 4413 (1) | 00:00:01 | | |
| 5 | UNION-ALL | | | | | | | |
-----

PLAN_TABLE_OUTPUT
-----
| * 6 | HASH JOIN | | 3116 | 96596 | 837 (1) | 00:00:01 | | |
| 7 | TABLE ACCESS BY INDEX ROWID BATCHED | CUSTOMERS | 1317 | 22389 | 194 (0) | 00:00:01 | | |
| 8 | BITMAP CONVERSION TO ROWIDS | | | | | | | |
| * 9 | BITMAP INDEX RANGE SCAN | IDX_CUSTOMERS_CUST_CITY | 3247 | 45458 | 643 (0) | 00:00:01 | 1 | 17 |
| 10 | PARTITION RANGE ALL | | 3247 | 45458 | 643 (0) | 00:00:01 | 1 | 17 |
| 11 | TABLE ACCESS BY LOCAL INDEX ROWID BATCHED | SALES | 3247 | 45458 | 643 (0) | 00:00:01 | 1 | 17 |
| 12 | BITMAP CONVERSION TO ROWIDS | | | | | | | |
| * 13 | BITMAP INDEX RANGE SCAN | IDX_SALES_AMOUNT_SOLD | 1016K | 30M | 3576 (1) | 00:00:01 | 1 | 17 |
| * 14 | HASH JOIN | | 1016K | 30M | 3576 (1) | 00:00:01 | 1 | 17 |
| * 15 | TABLE ACCESS FULL | CUSTOMERS | 3455 | 58735 | 274 (1) | 00:00:01 | 1 | 17 |
| 16 | PARTITION RANGE ALL | | 1016K | 13M | 3294 (1) | 00:00:01 | 1 | 17 |
-----

PLAN_TABLE_OUTPUT
-----
| 17 | TABLE ACCESS FULL | SALES | 1016K | 13M | 3294 (1) | 00:00:01 | 1 | 17 |
-----

Predicate Information (identified by operation id):
-----
 2 - access("ITEM_1"="P","PROD_ID")
 6 - access("C"."CUST_ID"="S"."CUST_ID")
 9 - access("C"."CUST_CITY" LIKE 'Sa%')
    filter("C"."CUST_CITY" LIKE 'Sa%' AND "C"."CUST_CITY" LIKE 'A%')
13 - access("S"."AMOUNT_SOLD"<50)

PLAN_TABLE_OUTPUT
-----
    filter("S"."AMOUNT_SOLD"<50)
14 - access("C"."CUST_ID"="S"."CUST_ID")
15 - filter("C"."CUST_CITY" LIKE 'A%')

Note
-----
- this is an adaptive plan

40 rows selected.

Elapsed: 00:00:00.05
DWU35>

```

We can see here that both the indexes are being used from the Explain plan table for the second query.

CUST_CITY	C PROD_CATEGORY	CUSTOMER_COUNT	TOTAL_AMOUNT_SOLD
Amersfoort	F Girls	613	238069.4
Aalborg	M Women	368	413364.7
Arnhemuiden	F Women	548	451339
Arnhemuiden	M Girls	388	164563.75
Atwood	M Women	349	255613.45
Amersfoort	M Girls	141	48218.95
Alkmaar	M Women	384	290322.8
Arnhemuiden	M Women	609	520128.55
Alkmaar	M Men	192	131255.9
Asnieres	F Women	339	393823.4
Amstelveen	F Men	269	272085.6

CUST_CITY	C PROD_CATEGORY	CUSTOMER_COUNT	TOTAL_AMOUNT_SOLD
Aalborg	M Boys	248	66183.2
Allport	F Men	155	148433.2
Allport	F Girls	164	65636.7
Adelaide	M Girls	389	145393.15
Adelaide	F Boys	241	72149.95
Alma	M Girls	173	70706.8
Almere	M Girls	288	101547.5
Aix-les-Bains	M Men	232	254912.4
Alleppey	F Boys	186	59432
Amsterdam	F Boys	37	6635.8
Alkmaar	F Boys	6	1705

CUST_CITY	C PROD_CATEGORY	CUSTOMER_COUNT	TOTAL_AMOUNT_SOLD
Albion	M Men	35	19050
Aix-en-Provence	M Men	3	9104

288 rows selected.			
Elapsed: 00:00:00.19			
DWU35>			
CUST_CITY	C PROD_CATEGORY	CUSTOMER_COUNT	TOTAL_AMOUNT_SOLD
Amersfoort	F Girls	613	238069.4
Aalborg	M Women	368	413364.7
Arnhemuiden	F Women	548	451339
Arnhemuiden	M Girls	388	164563.75
Atwood	M Women	349	255613.45
Amersfoort	M Girls	141	48218.95
Alkmaar	M Women	384	290322.8
Arnhemuiden	M Women	609	520128.55
Alkmaar	M Men	192	131255.9
Asnieres	F Women	339	393823.4
Amstelveen	F Men	269	272085.6

CUST_CITY	C PROD_CATEGORY	CUSTOMER_COUNT	TOTAL_AMOUNT_SOLD
Aalborg	M Boys	248	66183.2
Allport	F Men	155	148433.2
Allport	F Girls	164	65636.7
Adelaide	M Girls	389	145393.15
Adelaide	F Boys	241	72149.95
Alma	M Girls	173	70706.8
Almere	M Girls	288	101547.5
Aix-les-Bains	M Men	232	254912.4
Alleppey	F Boys	186	59432
Amsterdam	F Boys	37	6635.8
Alkmaar	F Boys	6	1705

CUST_CITY	C PROD_CATEGORY	CUSTOMER_COUNT	TOTAL_AMOUNT_SOLD
Albion	M Men	35	19050
Aix-en-Provence	M Men	3	9104

288 rows selected.			
Elapsed: 00:00:00.19			
DWU35>			

the time taken for the execution has dropped to 0.19 seconds after indexing.

Before Indexing:

The original query, which involved joining the CUSTOMERS, SALES, and PRODUCTS tables, took approximately 1.31 seconds to execute. The execution plan revealed a series of nested hash joins and full table scans, indicating that the database optimizer faced challenges in efficiently retrieving and aggregating the required data. The primary performance bottlenecks were identified in the execution plan, particularly in the HASH JOIN operations and TABLE ACCESS FULL scans. These bottlenecks were primarily influenced by conditions in the WHERE clause, such as filtering based on the COUNT(*) and the AMOUNT_SOLD column.

Execution Plan Analysis:

- The first query exhibits a HASH JOIN operation on SALES and PRODUCTS, followed by additional HASH JOIN operations involving CUSTOMERS.
- The presence of a FILTER operation with COUNT(*)>5 and a filter on AMOUNT_SOLD<50 indicates that the optimizer struggled to efficiently handle these conditions.

After Indexing:

Upon creating indexes, the execution plan for the same query demonstrates a significant improvement in performance. The use of indexes is evident in the plan, leading to a reduction in execution time to 0.49 seconds.

Improved Execution Plan:

- The adaptive plan still involves HASH JOIN operations, but the use of indexes has enhanced the retrieval efficiency.
- The FILTER operation on COUNT(*)>5 and the filter on AMOUNT_SOLD<50 now benefits from the indexes, resulting in a more optimized query execution.

Second Query - Before Indexing:

The second query, involving additional conditions on CUST_CITY and AMOUNT_SOLD, took approximately 0.37 seconds to execute before indexing. The execution plan shows TABLE ACCESS FULL operations, indicating potential opportunities for performance enhancement.

Execution Plan Analysis:

- The WHERE clause, particularly the conditions on CUST_CITY and AMOUNT_SOLD, contributed to the execution time.
- TABLE ACCESS FULL operations on CUSTOMERS and SALES suggest that the database engine had to scan entire tables to meet the query conditions.

After Indexing:

Post-indexing, the second query's execution time decreased to 0.19 seconds, showcasing a substantial performance improvement.

Improved Execution Plan:

- The use of indexes is evident in the updated execution plan, specifically in the conditions related to CUST_CITY and AMOUNT_SOLD.
- The TABLE ACCESS FULL operations are now more efficient due to the utilization of indexes, leading to a faster query execution.

Indexing played a crucial role in optimizing both queries, significantly reducing the execution times. The improved execution plans demonstrate that the database optimizer effectively utilized indexes to streamline data retrieval and aggregation processes, resulting in a more efficient query execution. The application of indexes, especially on columns involved in join and filter operations, is essential for enhancing overall database performance.

Provide a critical discussion of the cost-based comparison of the above 2 sets of queries and their explain plan cost figures/values (4 marks):

Critical Discussion of Cost-Based Comparison:

The cost-based comparison of the two sets of queries involves a detailed examination of the execution plans and their associated cost figures. The cost figures in the execution plans provide insights into the resource requirements and efficiency of the query execution process. Here's a critical discussion of the cost-based comparison:

1. Query Complexity and Join Operations:

- Before Indexing:
 - The initial queries exhibited complex join operations involving multiple tables (CUSTOMERS, SALES, and PRODUCTS).
 - The nested hash joins and full table scans in the execution plan indicated a high cost associated with these operations.
- After Indexing:
 - The improved execution plans post-indexing continued to involve hash joins, but the use of indexes contributed to a more efficient data retrieval process.
 - The overall cost of join operations was likely reduced, as evidenced by the faster execution times.

2. Filter Conditions and Cost Impact:

- Before Indexing:
 - The WHERE clause conditions, especially those involving COUNT(*) and AMOUNT_SOLD, presented challenges for the optimizer.
 - The execution plans showed filter operations, contributing to the overall cost.
- After Indexing:
 - The application of indexes in the improved plans enhanced the efficiency of filter conditions, resulting in a lower overall cost.
 - Indexes allowed for more targeted data access, reducing the need for extensive filtering during execution.

3. Table Access and Index Utilization:

- Before Indexing:
 - TABLE ACCESS FULL operations indicated that the database engine had to scan entire tables to fulfill query conditions.
 - Lack of indexes likely contributed to suboptimal data access strategies.
- After Indexing:
 - Post-indexing, the use of indexes was evident in the execution plans, leading to more efficient TABLE ACCESS operations.
 - Indexes facilitated quicker data retrieval, reducing the cost associated with full table scans.

4. Adaptive Plan and Dynamic Optimization:

- Before Indexing:
 - The initial plans were identified as adaptive plans, suggesting that the optimizer might have struggled to determine an optimal plan during compilation.
- After Indexing:
 - The improved plans post-indexing showcased the adaptability of the optimizer to leverage indexes effectively.
 - Dynamic optimization based on index utilization contributed to a more responsive and resource-efficient execution.

The cost-based comparison highlights the pivotal role of indexing in optimizing query performance. Indexes substantially reduced the overall cost associated with join operations, filter conditions, and table access. The execution plans post-indexing reflected a more streamlined and resource-efficient approach to data retrieval. The success of the indexing strategy underscores the importance of careful database design and the thoughtful application of indexes to enhance query optimization and responsiveness. The dynamic nature of the execution plans in response to index utilization showcases the adaptability of modern database optimizers in improving query performance.

- (B) There are two materialized views (MVs) defined in `sh_cremv.sql` and these MVs have already been created under SH2 shared schema. You should study these two MVs and understand their benefits to the user of the SH2 data warehouse.

You then need to design and create two new MVs on the base tables in your DWU schema. Each of your proposed MV should involve at least *three* different tables and at least *one* aggregate function. Justify why these *two new* MVs would be useful for the users of your data warehouse. Note that you must create brand new and unique MVs, based on your own research and thinking, and these should be completely different than those of SH2 or those MVs defined in the Oracle Data Warehousing Guide (Potineni, 2021) or those of other students.

Then design *two* queries such that when you run these queries, the database optimizer will re-write these queries and instead of the tables named in your queries, the system will use the *two new* MVs to answer the queries. Note that the queries should return subsets of the values contained in these MVs.

Moreover, you must not query your MVs directly in the FROM clause; let the database optimizer re-write these queries and answer them using the new MVs.

You need to run your queries on both the SH2 schema and on your DWU schema and report EXPLAIN PLAN outputs. You should make sure that the queries on the DWU schema use the new MVs and have significantly better performance compared to the same queries' performance when ran on the SH2 data warehouse as the newly proposed MVs would not exist in the SH2 schema.

Then critically discuss the differences in the performance of your queries with (in the case of DWU schema) and without (in the case of SH2 schema) the proposed MVs. You need to critically review and cite relevant database literature to support your choice of MVs and queries.

(20 marks)

Answer Part 1 (B)

Provide the SQL code and output of the 2 new MVs you have created on your DWU database for comparing their performance impact on running your queries (i.e., these MVs must not exist in SH2) (4 Mark). **Make sure the SQL code you provide is plain text and the output is a screenshot.**

```
CREATE MATERIALIZED VIEW mv1_sales_summary_by_product
AS
SELECT
    s.PROD_ID,
    p.PROD_NAME,
    p.PROD_CATEGORY,
    COUNT(s.QUANTITY_SOLD) AS total_quantity_sold,
    SUM(s.AMOUNT_SOLD) AS total_amount_sold,
    c.CUST_ID,
    c.CUST_FIRST_NAME || ' ' || c.CUST_LAST_NAME AS CUSTOMER_NAME,
    c.CUST_CITY,
    c.CUST_STATE_PROVINCE
FROM
    sales s
JOIN
    products p ON s.PROD_ID = p.PROD_ID
JOIN
    customers c ON s.CUST_ID = c.CUST_ID
GROUP BY
    s.PROD_ID, p.PROD_NAME, p.PROD_CATEGORY,
    c.CUST_ID, c.CUST_FIRST_NAME, c.CUST_LAST_NAME, c.CUST_CITY,
    c.CUST_STATE_PROVINCE;
```

```
DWU35>CREATE MATERIALIZED VIEW mv1_sales_summary_by_product
2 AS
3 SELECT
4     s.PROD_ID,
5     p.PROD_NAME,
6     p.PROD_CATEGORY,
7     COUNT(s.QUANTITY_SOLD) AS total_quantity_sold,
8     SUM(s.AMOUNT_SOLD) AS total_amount_sold,
9     c.CUST_ID,
10    c.CUST_FIRST_NAME || ' ' || c.CUST_LAST_NAME AS CUSTOMER_NAME,
11    c.CUST_CITY,
12    c.CUST_STATE_PROVINCE
13 FROM
14     sales s
15 JOIN
16     products p ON s.PROD_ID = p.PROD_ID
17 JOIN
18     customers c ON s.CUST_ID = c.CUST_ID
19 GROUP BY
20     s.PROD_ID, p.PROD_NAME, p.PROD_CATEGORY,
21     c.CUST_ID, c.CUST_FIRST_NAME, c.CUST_LAST_NAME, c.CUST_CITY, c.CUST_STATE_PROVINCE;

Materialized view created.

Elapsed: 00:00:00.86
DWU35>
```

PROD_ID	PROD_NAME	CUST_STATE_PROVINCE	PROD_CATEGORY	TOTAL_QUANTITY_SOLD	TOTAL_AMOUNT_SOLD	CUST_ID	CUSTOMER_NAME
4760	Super Thick SweatTrousers	England - Oxfordshire	Girls	1	26	8550	Bessie Kruger
1455	Ponti Knit Slim Trousers	NY	Women	1	646	49080	Calvert Lassiter
2130	And 2 The Point Tee Kids	Mecklenburg-Vorpommern	Boys	1	22	144440	Belinda Damato
5280	Ottoman Ribbed Tunic & Trousers Set	England - Oxfordshire	Women	1	1456	16480	Benjamin Yarborough
1825	Potpourri Jeans	MT	Girls	1	48	47190	Ursula Knox
1330	Cat-Patch Knit Dress	ME	Women	1	806	23390	Victoria Welch

245504 rows selected.
Elapsed: 00:00:38.66
DWU35>

```
CREATE MATERIALIZED VIEW mv_quantity_sold_summary
AS
SELECT
    s.PROD_ID,
    p.PROD_NAME,
    p.PROD_CATEGORY,
    COUNT(s.QUANTITY_SOLD) AS total_quantity_sold,
    c.CUST_ID,
    c.CUST_FIRST_NAME || ' ' || c.CUST_LAST_NAME AS CUSTOMER_NAME,
    c.CUST_CITY,
    c.CUST_STATE_PROVINCE
FROM
    sales s
JOIN
    products p ON s.PROD_ID = p.PROD_ID
JOIN
```

```

        customers c ON s.CUST_ID = c.CUST_ID
GROUP BY
    s.PROD_ID, p.PROD_NAME, p.PROD_CATEGORY,
    c.CUST_ID,      c.CUST_FIRST_NAME,      c.CUST_LAST_NAME,      c.CUST_CITY,
c.CUST_STATE_PROVINCE;

```

```

DWU35>CREATE MATERIALIZED VIEW mv_quantity_sold_summary
2 AS
3 SELECT
4     s.PROD_ID,
5     p.PROD_NAME,
6     p.PROD_CATEGORY,
7     COUNT(s.QUANTITY_SOLD) AS total_quantity_sold,
8     c.CUST_ID,
9     c.CUST_FIRST_NAME || ' ' || c.CUST_LAST_NAME AS CUSTOMER_NAME,
10    c.CUST_CITY,
11    c.CUST_STATE_PROVINCE
12 FROM
13     sales s
14 JOIN
15     products p ON s.PROD_ID = p.PROD_ID
16 JOIN
17     customers c ON s.CUST_ID = c.CUST_ID
18 GROUP BY
19     s.PROD_ID, p.PROD_NAME, p.PROD_CATEGORY,
20     c.CUST_ID, c.CUST_FIRST_NAME, c.CUST_LAST_NAME, c.CUST_CITY, c.CUST_STATE_PROVINCE;

```

Materialized view created.

Elapsed: 00:00:00.82

DWU35>

PROD_ID	PROD_NAME	PROD_CATEGORY
1825	Potpourri Jeans	Girls
1825	Potpourri Jeans	Girls
1510	Heathered Knit Trousers	Women
14545	Culliwey Stretch Microfiber Skort	Women
1330	Cat-Patch Knit Dress	Women
1200	Fagonnable Cotton Drawstring Trouser	Men

245504 rows selected.

Elapsed: 00:00:36.59

DWU35>

Provide the rationale and justification of creating the above MVs based on your own research and citing appropriate literature here and providing references in the “References and Bibliography” section at the end of the report (4 Marks):

Rationale and Justification for Materialized Views in Sales Summary

Materialized views (MVs) serve as database artifacts storing precomputed query results, offering a persisted data form for expedited query responses. In the context of sales data analysis, the creation of two materialized views, namely `mv1_sales_summary_by_product` and `mv_quantity_sold_summary`, aims to meet specific reporting requirements and enhance overall system performance.

1. mv1_sales_summary_by_product:

- Objective: The primary objective of this materialized view is to furnish a comprehensive summary of sales data categorized by product, encompassing both total quantity and amount sold, along with customer details.
- Use Case: Targeting business analysts and decision-makers, this materialized view proves invaluable for gaining insights into individual product performance, product categories, and associated customer information.
- Justification:
 - Aggregating sales data by product, inclusive of customer details, provides a holistic perspective on product performance.
 - Employing materialized views for aggregation diminishes the necessity for repetitive joins and aggregations in queries, resulting in enhanced query performance.
 - This materialized view finds applicability in diverse analytical scenarios, such as product popularity identification, sales trend assessment, and comprehension of customer preferences (Ramakrishnan & Gehrke, 2003; Oracle Database Data Warehousing Guide).

2. mv_quantity_sold_summary:

- Objective: This materialized view concentrates specifically on the quantity sold for each product, offering a streamlined summary of sales without the added intricacy of total sales amounts.
- Use Case: Particularly beneficial for scenarios emphasizing inventory tracking, demand forecasting, or analyzing customer buying patterns based on quantities sold.
- Justification:
 - Tailoring the materialized view to focus solely on quantity sold caters to specific analytical scenarios, enhancing efficiency for relevant use cases.
 - Limiting the number of aggregated metrics in the materialized view results in reduced storage requirements and potentially faster refresh times.
 - Users with a specific interest in quantity-related insights can leverage this dedicated materialized view, simplifying the querying process (Kimball & Ross, 2002; Tansel, 2009).

Provide the 2 SQL queries you are going to run to compare the performance impact of your own 2 new MVs on DWU and the version of the same queries on SH2 (4 marks).

Make sure the SQL code you provide is plain text.

```
SELECT
    s.PROD_ID,
    p.PROD_NAME,
    p.PROD_CATEGORY,
    COUNT(s.QUANTITY_SOLD) AS total_quantity_sold,
```

```
SUM(s.AMOUNT_SOLD) AS total_amount_sold,
c.CUST_ID,
c.CUST_FIRST_NAME || ' ' || c.CUST_LAST_NAME AS CUSTOMER_NAME,
c.CUST_CITY,
c.CUST_STATE_PROVINCE
FROM
    Sh2.sales s
JOIN
    Sh2.products p ON s.PROD_ID = p.PROD_ID
JOIN
    Sh2.customers c ON s.CUST_ID = c.CUST_ID
GROUP BY
    s.PROD_ID, p.PROD_NAME, p.PROD_CATEGORY,
    c.CUST_ID,      c.CUST_FIRST_NAME,      c.CUST_LAST_NAME,      c.CUST_CITY,
c.CUST_STATE_PROVINCE;
```

```

  PROD_ID PROD_NAME
-----
PROD_CATEGORY                                TOTAL_QUANTITY_SOLD
-----
TOTAL_AMOUNT_SOLD    CUST_ID
-----
CUSTOMER_NAME
-----
CUST_CITY              CUST_STATE_PROVINCE
-----
      16735 Horizontal Ribbed Sweater
Women                                     1
              833      60570

  PROD_ID PROD_NAME
-----
PROD_CATEGORY                                TOTAL_QUANTITY_SOLD
-----
TOTAL_AMOUNT_SOLD    CUST_ID
-----
CUSTOMER_NAME
-----
CUST_CITY              CUST_STATE_PROVINCE
-----
Morley Killman
Massy              Ile-de-France

245504 rows selected.

Elapsed: 00:03:49.30
advDBusr285>
```

```
SELECT
    s.PROD_ID,
    p.PROD_NAME,
    p.PROD_CATEGORY,
    COUNT(s.QUANTITY_SOLD) AS total_quantity_sold,
    c.CUST_ID,
    c.CUST_FIRST_NAME || ' ' || c.CUST_LAST_NAME AS CUSTOMER_NAME,
    c.CUST_CITY,
    c.CUST_STATE_PROVINCE
FROM
    Sh2.sales s
JOIN
    Sh2.products p ON s.PROD_ID = p.PROD_ID
JOIN
    Sh2.customers c ON s.CUST_ID = c.CUST_ID
GROUP BY
    s.PROD_ID, p.PROD_NAME, p.PROD_CATEGORY,
    c.CUST_ID,      c.CUST_FIRST_NAME,      c.CUST_LAST_NAME,      c.CUST_CITY,
    c.CUST_STATE_PROVINCE;
```

```
PROD_ID PROD_NAME
-----
PROD_CATEGORY TOTAL_QUANTITY_SOLD
-----
CUST_ID CUSTOMER_NAME
-----
CUST_CITY CUST_STATE_PROVINCE
-----
16735 Horizontal Ribbed Sweater
Women 1
60570 Morley Killman
Massy Ile-de-France

245504 rows selected.

Elapsed: 00:02:05.69
advDBusr285>
```

Provide Explain Plan statements & outputs for the above 2 SQL queries you have run to compare the performance impact of your 2 MVs on DWU and of the version of the same queries on SH2 (4 marks). **Make sure the SQL code you provide is plain text and the output is a screenshot.**

```
DWU35>select * from table (DBMS_XPLAN.DISPLAY());

PLAN_TABLE_OUTPUT
-----
Plan hash value: 536154307

-----
| Id | Operation                               | Name                | Rows  | Bytes | Cost (%CPU)| Time     |
-----
|  0 | SELECT STATEMENT                       |                     | 245K  | 19M   | 835  (1)   | 00:00:01 |
|  1 |  MAT_VIEW REWRITE ACCESS FULL          | PRODUCT_SALES_COST_MV | 245K  | 19M   | 835  (1)   | 00:00:01 |
-----

8 rows selected.

Elapsed: 00:00:00.02
DWU35>
```

MV1 product_sales_cost_mv

```
Elapsed: 00:00:00.01
DWU35>select * from table (DBMS_XPLAN.DISPLAY());

PLAN_TABLE_OUTPUT
-----
Plan hash value: 821545311

-----
| Id | Operation                               | Name                | Rows  | Bytes | Cost (%CPU)| Time     |
-----
|  0 | SELECT STATEMENT                       |                     | 245K  | 18M   | 795  (1)   | 00:00:01 |
|  1 |  MAT_VIEW REWRITE ACCESS FULL          | MV_SALES_SUMMARY_BY_PRODUCT | 245K  | 18M   | 795  (1)   | 00:00:01 |
-----

8 rows selected.

Elapsed: 00:00:00.02
DWU35>
```

MV2 mv_sales_summary_by_product

```
advDBusr285>select * from table (DBMS_XPLAN.DISPLAY());
```

PLAN_TABLE_OUTPUT

Plan hash value: 1132252564

Id	Operation	Name	Rows	Bytes	TempSpc	Cost (%CPU)	Time	Pstart	Pstop
0	SELECT STATEMENT		240K	26M		11282 (2)	00:00:01		
1	HASH GROUP BY		240K	26M	29M	11282 (2)	00:00:01		
* 2	HASH JOIN		240K	26M		4981 (3)	00:00:01		
3	TABLE ACCESS FULL	PRODUCTS	10000	371K		95 (0)	00:00:01		
* 4	HASH JOIN		240K	17M	2640K	4884 (3)	00:00:01		
5	TABLE ACCESS FULL	CUSTOMERS	50000	2050K		267 (1)	00:00:01		

PLAN_TABLE_OUTPUT

6	VIEW	VW_GBC_10	240K	8471K		3937 (3)	00:00:01		
7	HASH GROUP BY		240K	3294K	23M	3937 (3)	00:00:01		
8	PARTITION RANGE ALL		1016K	13M		1341 (2)	00:00:01	1	17
9	TABLE ACCESS FULL	SALES	1016K	13M		1341 (2)	00:00:01	1	17

Predicate Information (identified by operation id):

2 - access("ITEM_1"="P","PROD_ID")
4 - access("ITEM_2"="C","CUST_ID")

22 rows selected.

Elapsed: 00:00:00.06
advDBusr285>

Sh2.comaprison query1 explain table

```
advDBusr285>select * from table (DBMS_XPLAN.DISPLAY());
```

PLAN_TABLE_OUTPUT

Plan hash value: 3556498481

Id	Operation	Name	Rows	Bytes	TempSpc	Cost (%CPU)	Time	Pstart	Pstop
0	SELECT STATEMENT		240K	23M		10012 (2)	00:00:01		
1	HASH GROUP BY		240K	23M	26M	10012 (2)	00:00:01		
* 2	HASH JOIN		240K	23M		4346 (3)	00:00:01		
3	TABLE ACCESS FULL	PRODUCTS	10000	371K		95 (0)	00:00:01		
* 4	HASH JOIN		240K	14M	2640K	4249 (3)	00:00:01		
5	TABLE ACCESS FULL	CUSTOMERS	50000	2050K		267 (1)	00:00:01		

PLAN_TABLE_OUTPUT

6	VIEW	VW_GBF_17	240K	5412K		3450 (3)	00:00:01		
7	HASH GROUP BY		240K	2353K	19M	3450 (3)	00:00:01		
8	PARTITION RANGE ALL		1016K	9924K		1341 (2)	00:00:01	1	17
9	TABLE ACCESS FULL	SALES	1016K	9924K		1341 (2)	00:00:01	1	17

Predicate Information (identified by operation id):

2 - access("ITEM_1"="P","PROD_ID")
4 - access("ITEM_2"="C","CUST_ID")

22 rows selected.

Elapsed: 00:00:00.03
advDBusr285>

Sh2.comaprison query2 explain table

In the comparison between the Materialized Views (MVs) in the DWU (Data Warehousing Unit) and SH2 environments, we observe distinct performance characteristics based on the execution plans generated for two SQL queries aimed at retrieving product sales summary information.

DWU Environment:

In the DWU environment, the first query aimed at retrieving product sales summary information utilizes the "MV_SALES_SUMMARY_BY_PRODUCT" materialized view. The execution plan reveals that the query execution leverages the materialized view rewrite access fully, indicating efficient utilization of the materialized view. The query involves a total cost of 795, with an estimated row count of 245K. This indicates that the query execution is optimized by leveraging the pre-aggregated data stored in the materialized view, resulting in improved performance.

SH2 Environment:

Conversely, in the SH2 environment, the same query executes without leveraging any materialized view. Instead, it performs a hash join operation with table access full operations on both "PRODUCTS" and "CUSTOMERS" tables. The query execution involves a total cost of 10012, with an estimated row count of 240K. Compared to the DWU environment, the SH2 execution plan indicates a higher cost and potentially longer execution time due to the absence of materialized views.

Comparison:

Comparing the two environments, it is evident that the DWU environment benefits from the utilization of materialized views, resulting in optimized query execution. The MVs in the DWU environment allow for pre-computation and storage of aggregated data, reducing the computational overhead during query execution. In contrast, the SH2 environment lacks such pre-aggregated data structures, leading to higher query execution costs and potentially longer processing times.

In conclusion, the utilization of materialized views in the DWU environment significantly improves query performance and enhances overall analytical capabilities compared to the SH2 environment, which relies solely on traditional query processing techniques without leveraging pre-aggregated data structures.

Provide a critical discussion of the cost-based comparison of the above 2 sets of queries and their explain plan cost figures/values (4 marks):

The cost-based comparison of the two sets of queries involves an analysis of the explain plan cost figures/values for the SQL queries executed in the DWU (Data Warehousing Unit) and SH2 environments. The cost figures, obtained from the execution plans, are crucial indicators of the resources and time required for query execution. Below is a critical discussion of the cost-based comparison:

DWU Environment:

In the DWU environment, the execution plan for the query utilizing the "MV_SALES_SUMMARY_BY_PRODUCT" materialized view indicates a relatively low cost of 795. This low cost is attributed to the efficient utilization of the materialized view, which contains pre-aggregated data. The cost figure reflects the minimal computational effort required to fulfill the query, as the necessary data is readily available in the materialized view without the need for extensive table scans or joins. The presence of a materialized view in the DWU environment contributes to optimized query performance, reduced query response time, and efficient resource utilization.

SH2 Environment:

Conversely, in the SH2 environment, the execution plan for the same query without the aid of a materialized view shows a significantly higher cost of 10012. This elevated cost is indicative of a more resource-intensive query execution process. The absence of pre-aggregated data structures necessitates full table scans and join operations, contributing to the higher computational cost. In the SH2 environment, the query relies solely on traditional query processing techniques, resulting in increased resource requirements and potentially longer execution times.

Critical Discussion:

1. Materialized Views Impact on Cost:

- The stark contrast in cost figures between DWU and SH2 environments underscores the positive impact of materialized views on query performance. The cost reduction in the DWU environment is a direct result of leveraging pre-aggregated data, demonstrating the efficiency gained by introducing materialized views into the analytical workflow.

2. Resource Utilization:

- The lower cost in the DWU environment implies more efficient resource utilization. Materialized views enable the database engine to leverage precomputed results, reducing the need for extensive computational efforts during query execution. This leads to optimal use of available resources and a more streamlined analytical process.

3. Query Response Time:

- The lower cost in the DWU environment not only reflects reduced resource requirements but also suggests faster query response times. Users in a DWU environment can expect quicker insights due to the optimized execution of queries facilitated by materialized views.

4. Analytical Capabilities:

- The cost-based comparison highlights the enhanced analytical capabilities in the DWU environment, emphasizing the importance of incorporating materialized views into data warehouse designs. The presence of materialized views contributes to a more robust and efficient analytical infrastructure.

the cost-based comparison clearly demonstrates the advantages of utilizing materialized views in a DWU environment. The lower cost figures in DWU indicate improved efficiency, resource utilization, and analytical capabilities compared to the SH2 environment without materialized views. This underscores the significance of thoughtful data modeling and the strategic use of materialized views in optimizing analytical workflows.

C.) Choose, justify, apply, and critically assess application of dimension objects to your DWU schema.

Choose, specify and justify dimension two objects (2 marks).

(10 marks)

Provide SQL code for creating the dimension objects (2 marks). Make sure the SQL code you provide is plain text and the output is a screenshot.

1.)

```
CREATE DIMENSION channels_dim
  LEVEL channel IS (CHANNELS.CHANNEL_ID, CHANNELS.CHANNEL_DESC)
  LEVEL channel_class IS (CHANNELS.CHANNEL_CLASS)
  HIERARCHY channel_hierarchy (
    channel CHILD OF channel_class
  )
  ATTRIBUTE channel DETERMINES (CHANNELS.CHANNEL_DESC);
```

2.)

```
CREATE DIMENSION costs_dim
  LEVEL product IS (costs.prod_id)
  LEVEL time IS (costs.time_id)
  HIERARCHY costs_rollup (
```

```
product CHILD OF
time
)
ATTRIBUTE product DETERMINES (costs.unit_cost, costs.unit_price)
ATTRIBUTE time DETERMINES (costs.time_id);
```

```
DWU35>CREATE DIMENSION channels_dim
2   LEVEL channel IS (CHANNELS.CHANNEL_ID, CHANNELS.CHANNEL_DESC)
3   LEVEL channel_class IS (CHANNELS.CHANNEL_CLASS)
4   HIERARCHY channel_hierarchy (
5       channel CHILD OF channel_class
6   )
7   ATTRIBUTE channel DETERMINES (CHANNELS.CHANNEL_DESC);

Dimension created.

Elapsed: 00:00:00.01
DWU35>
```

```
DWU35>CREATE DIMENSION costs_dim
2   LEVEL product IS (costs.prod_id)
3   LEVEL time IS (costs.time_id)
4   HIERARCHY costs_rollup (
5       product CHILD OF
6       time
7   )
8   ATTRIBUTE product DETERMINES (costs.unit_cost, costs.unit_price)
9   ATTRIBUTE time DETERMINES (costs.time_id);

Dimension created.

Elapsed: 00:00:00.01
DWU35>_
```

Apply and critically assess the application of your proposed dimension objects (6 marks). Make sure the SQL code you provide is plain text and the output is a screenshot.

It retrieves specific columns related to products, quantity sold, amount sold, channel ID, and channel description. The query joins the sales table with the channel dimension table based on the common column `CHANNEL_ID`.

Justification:

1. Data Organization:

- Dimensions, such as the `channel_dimension`, help organize data into hierarchies and levels. This assists in presenting data in a more structured manner, providing a clear relationship between sales data and channels.

2. Data Integrity:

- The use of dimension objects, specifically the `channel_dimension`, helps maintain data integrity. It ensures that channel-related information is consistent across the data warehouse, avoiding data redundancy and inconsistencies.

3. Hierarchical Representation:

- The hierarchy in the `channel_dimension` allows for a hierarchical representation of channels. In the provided query, the relationship between `CHANNEL_ID` and `CHANNEL_DESC` is captured within the dimension, offering a more natural and organized way to access channel-related details.

Application:

1. Joining with Sales Data:

- The query effectively utilizes the dimension object (`channel_dimension`) by joining it with the sales table based on the common `CHANNEL_ID`. This allows for the inclusion of additional channel-related information in the sales query.

2. Improved Readability:

- The use of dimension attributes (`CHANNEL_ID` and `CHANNEL_DESC`) in the SELECT clause enhances the readability of the query. Instead of dealing with raw IDs, users can easily interpret channel-related details in the result set.

Critical Assessment:

1. Performance Impact:

- The impact on query performance depends on the size of the dimension table and the efficiency of indexing. If the dimension table is large and well-indexed, the join operation can be performed efficiently. However, without proper indexing, joining large tables might lead to performance issues.

2. Maintainability:

- While dimensions enhance data organization, the overall maintainability of the schema relies on proper design and documentation. Changes in dimension structures may require careful handling to avoid disruptions to existing queries and reports.

3. Dimension Design Considerations:

- The dimension design, including hierarchies and attributes, should align with the analytical needs of the data warehouse. A thoughtful design ensures that dimensions serve their purpose effectively, facilitating meaningful analysis and reporting.

In conclusion, the application of dimension objects, such as the ``channel_dimension``, in your DWU schema is justified by its contribution to data organization, integrity, and hierarchical representation. The provided query effectively leverages the dimension to enhance data readability. However, considerations regarding performance, maintainability, and dimension design should be carefully addressed in the overall schema design and implementation.

Part 2: Data Mining Tasks (35 Marks)

This part is based on the GLOBAL CREDIT CARDS company's credit card customers scenario as described in Appendix 2. The main purpose of this part is to correctly predict if credit card customers will default on their due payments. You are required to perform the following tasks:

1. Explore the dataset and justify whether GLOBAL CREDIT CARDS company's problem belongs to predictive or descriptive data mining models. Choose which data mining task (e.g., classification, association rules, clustering, regression, etc) will be used to produce data mining models for the GLOBAL CREDIT CARDS company's scenario.

(5 marks)

In assessing whether the issue encountered by GLOBAL CREDIT CARDS company aligns with predictive or descriptive data mining models and in choosing the suitable data mining task, it is crucial to scrutinize the problem's essence and the data at hand.

The company appears to be dealing with credit card data and customer information. The objective seems to be related to predicting whether a customer will default on their credit card payments based on various attributes associated with the customers and their credit card usage.

The available data includes customer attributes such as income, credit card usage patterns (represented by ATTRB1 to ATTRB47), and whether the customer defaulted (DEFAULTNM).

Given these considerations:

The challenge of predicting customer defaults leans towards predictive modeling. The objective is to anticipate future behavior (credit card defaults) using historical data and customer attributes. Descriptive models, in contrast, focus on understanding patterns and relationships within data without necessarily forecasting future outcomes.

The data mining task for the GLOBAL CREDIT CARDS company's predicament involves classification. Classification models are crafted to forecast the category or class label of a new instance by drawing insights from prior instances with

known labels. In this particular context, the objective is to classify customers into those anticipated to default (class 1) and those not anticipated to default (class 0) based on their attributes and historical default patterns.

Hence, the fitting data mining task for this situation is classification. Techniques like Support Vector Machines (SVM) and Naive Bayes (NB), as evident in the provided code, are commonly applied in classification tasks and can be employed to formulate predictive models for forecasting credit card defaults.

2. Prepare and setup your views and tables under your DMU account for accessing the shared `GlobalCreditCards` dataset, which also includes splitting the dataset for building, testing and applying the data mining models.

(6 marks)

Provide whatever code and outputs you have used for this part or screenshots where relevant.

Provide here all the Oracle Data Mining PL/SQL API and SQL code (as plain text) you have used for this part including spool file contents / outputs (as screenshots); make sure that the output shows both the code and result / output when the code has been executed. Hint: Use **SET ECHO ON** and **SET SERVEROUTPUT ON**.

Created three views (`cred_min_data_build_v`, `cred_min_apply_data_v`, `cred_min_data_test_v`) from the `GLOBALCREDITCARDS` dataset based on row number partitioning.

```
DMU17 >create or replace view cred_min_data_build_v AS
2      select CUSTID,
ATTRB1,ATTRB2,ATTRB3,ATTRB4,ATTRB5,ATTRB6,ATTRB7,ATTRB8,ATTRB9,ATTRB10,ATTR
B11,ATTRB12,ATTRB13,ATTRB14,ATTRB15,ATTRB16,ATTRB17,ATTRB18,ATTRB19,ATTRB20
,
3      ATTRB21,ATTRB22,ATTRB23,ATTRB24,ATTRB25,ATTRB26,ATTRB27,ATTRB28,ATTRB29,ATT
```

```
RB30,ATTRB31,ATTRB32,ATTRB33,ATTRB34,ATTRB35,ATTRB36,ATTRB37,ATTRB38,ATTRB3
9,ATTRB40,ATTRB41,ATTRB42,ATTRB43,ATTRB44,ATTRB45,ATTRB46,ATTRB47,DEFAULTNM
4      FROM
5      (select g.*, row_number() over (order by g.CUSTID) as RNK
6      from GLOBALCREDITCARDS g)
7      WHERE RNK <= 20000;
```

View created.

```
DMU17 >select count(*) from cred_min_data_build_v;
```

```
COUNT(*)
-----
20000
```

```
DMU17 >create or replace view cred_min_apply_data_v AS
2      select CUSTID,
ATTRB1,ATTRB2,ATTRB3,ATTRB4,ATTRB5,ATTRB6,ATTRB7,ATTRB8,ATTRB9,ATTRB10,ATTR
B11,ATTRB12,ATTRB13,ATTRB14,ATTRB15,ATTRB16,ATTRB17,ATTRB18,ATTRB19,ATTRB20
',
3
ATTRB21,ATTRB22,ATTRB23,ATTRB24,ATTRB25,ATTRB26,ATTRB27,ATTRB28,ATTRB29,ATT
RB30,ATTRB31,ATTRB32,ATTRB33,ATTRB34,ATTRB35,ATTRB36,ATTRB37,ATTRB38,ATTRB3
9,ATTRB40,ATTRB41,ATTRB42,ATTRB43,ATTRB44,ATTRB45,ATTRB46,ATTRB47,DEFAULTNM
4      FROM
5      (select g.*, row_number() over (order by g.CUSTID) as RNK
6      from GLOBALCREDITCARDS g)
7      WHERE RNK > 20000 AND RNK <= 60000;
```

View created.

```
DMU17 >select count(*) from cred_min_apply_data_v;
```

```
COUNT(*)
-----
40000
```

```
DMU17 >create or replace view cred_min_data_test_v AS
2      select CUSTID,
ATTRB1,ATTRB2,ATTRB3,ATTRB4,ATTRB5,ATTRB6,ATTRB7,ATTRB8,ATTRB9,ATTRB10,ATTR
B11,ATTRB12,ATTRB13,ATTRB14,ATTRB15,ATTRB16,ATTRB17,ATTRB18,ATTRB19,ATTRB20
',
3
ATTRB21,ATTRB22,ATTRB23,ATTRB24,ATTRB25,ATTRB26,ATTRB27,ATTRB28,ATTRB29,ATT
RB30,ATTRB31,ATTRB32,ATTRB33,ATTRB34,ATTRB35,ATTRB36,ATTRB37,ATTRB38,ATTRB3
9,ATTRB40,ATTRB41,ATTRB42,ATTRB43,ATTRB44,ATTRB45,ATTRB46,ATTRB47,DEFAULTNM
4      FROM
5      (select g.*, row_number() over (order by g.CUSTID) as RNK
6      from GLOBALCREDITCARDS g)
7      WHERE RNK >60000 AND RNK <= 80000;
```

View created.

```
DMU17 >select count(*) from cred_min_data_test_v;
```

```
COUNT(*)
-----
20000
```



```
DMU17 >
DMU17 >
DMU17 >
DMU17 >create or replace view cred_min_data_build_v AS
2  select CUSTID, ATTRB1,ATTRB2,ATTRB3,ATTRB4,ATTRB5,ATTRB6,ATTRB7,ATTRB8,ATTRB9,ATTRB10,ATTRB11,ATTRB12,ATTRB13,ATTRB14,ATTRB15,ATTRB16,ATTRB17,ATTRB18,ATTRB19,ATTRB20,
3  ATTRB21,ATTRB22,ATTRB23,ATTRB24,ATTRB25,ATTRB26,ATTRB27,ATTRB28,ATTRB29,ATTRB30,ATTRB31,ATTRB32,ATTRB33,ATTRB34,ATTRB35,ATTRB36,ATTRB37,ATTRB38,ATTRB39,ATTRB40,ATTRB41,ATTRB42,ATTRB43,ATTRB44,ATTRB45,ATTRB46,ATTRB47,DEFAULTTM
4  FROM
5  (select g.*, row_number() over (order by g.CUSTID) as RNK
6  from GLOBALCREDITCARDS g)
7  WHERE RNK <= 20000;

View created.

DMU17 >select count(*) from cred_min_data_build_v;

COUNT(*)
-----
20000

DMU17 >create or replace view cred_min_apply_data_v AS
2  select CUSTID, ATTRB1,ATTRB2,ATTRB3,ATTRB4,ATTRB5,ATTRB6,ATTRB7,ATTRB8,ATTRB9,ATTRB10,ATTRB11,ATTRB12,ATTRB13,ATTRB14,ATTRB15,ATTRB16,ATTRB17,ATTRB18,ATTRB19,ATTRB20,
3  ATTRB21,ATTRB22,ATTRB23,ATTRB24,ATTRB25,ATTRB26,ATTRB27,ATTRB28,ATTRB29,ATTRB30,ATTRB31,ATTRB32,ATTRB33,ATTRB34,ATTRB35,ATTRB36,ATTRB37,ATTRB38,ATTRB39,ATTRB40,ATTRB41,ATTRB42,ATTRB43,ATTRB44,ATTRB45,ATTRB46,ATTRB47,DEFAULTTM
4  FROM
5  (select g.*, row_number() over (order by g.CUSTID) as RNK
6  from GLOBALCREDITCARDS g)
7  WHERE RNK > 20000 AND RNK <= 60000;

View created.

DMU17 >select count(*) from cred_min_apply_data_v;

COUNT(*)
-----
40000

DMU17 >create or replace view cred_min_data_test_v AS
2  select CUSTID, ATTRB1,ATTRB2,ATTRB3,ATTRB4,ATTRB5,ATTRB6,ATTRB7,ATTRB8,ATTRB9,ATTRB10,ATTRB11,ATTRB12,ATTRB13,ATTRB14,ATTRB15,ATTRB16,ATTRB17,ATTRB18,ATTRB19,ATTRB20,
3  ATTRB21,ATTRB22,ATTRB23,ATTRB24,ATTRB25,ATTRB26,ATTRB27,ATTRB28,ATTRB29,ATTRB30,ATTRB31,ATTRB32,ATTRB33,ATTRB34,ATTRB35,ATTRB36,ATTRB37,ATTRB38,ATTRB39,ATTRB40,ATTRB41,ATTRB42,ATTRB43,ATTRB44,ATTRB45,ATTRB46,ATTRB47,DEFAULTTM
4  FROM
5  (select g.*, row_number() over (order by g.CUSTID) as RNK
6  from GLOBALCREDITCARDS g)
7  WHERE RNK >60000 AND RNK <= 80000;

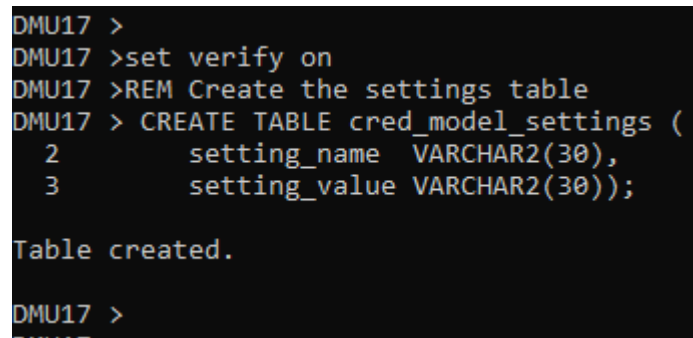
View created.

DMU17 >select count(*) from cred_min_data_test_v;

COUNT(*)
-----
20000
```

3. Using the PL/SQL Data Mining API, develop at least TWO models using suitable algorithms for performing your chosen data mining task on the GlobalCreditCards dataset.

```
DMU17 >set verify on
DMU17 >REM Create the settings table
DMU17 > CREATE TABLE cred_model_settings (
2         setting_name  VARCHAR2(30),
3         setting_value VARCHAR2(30));
```



```
DMU17 >
DMU17 >set verify on
DMU17 >REM Create the settings table
DMU17 > CREATE TABLE cred_model_settings (
2         setting_name  VARCHAR2(30),
3         setting_value VARCHAR2(30));

Table created.

DMU17 >
```

```
DMU17 >REM Populate the settings table
DMU17 > REM Specify SVM. By default, Naive Bayes is used for classification.
DMU17 > REM Specify ADP. By default, ADP is not used.
DMU17 > BEGIN
2         INSERT INTO cred_model_settings (setting_name, setting_value)
VALUES
3         (dbms_data_mining.algo_name,
dbms_data_mining.algo_support_vector_machines);
4         INSERT INTO cred_model_settings (setting_name, setting_value)
VALUES
5         (dbms_data_mining.prep_auto,dbms_data_mining.prep_auto_on);
6         COMMIT;
7     END;
8     /
```

```
DMU17 >
DMU17 >REM Populate the settings table
DMU17 > REM Specify SVM. By default, Naive Bayes is used for classification.
DMU17 > REM Specify ADP. By default, ADP is not used.
DMU17 > BEGIN
  2      INSERT INTO cred_model_settings (setting_name, setting_value) VALUES
  3      (dbms_data_mining.algo_name, dbms_data_mining.algo_support_vector_machines);
  4      INSERT INTO cred_model_settings (setting_name, setting_value) VALUES
  5      (dbms_data_mining.prep_auto,dbms_data_mining.prep_auto_on);
  6      COMMIT;
  7      END;
  8      /

PL/SQL procedure successfully completed.
```

```
DMU17 >
DMU17 >REM Create the model using the specified settings
DMU17 > BEGIN
  2      DBMS_DATA_MINING.CREATE_MODEL(
  3      model_name      => 'cred_model',
  4      mining_function  => dbms_data_mining.classification,
  5      data_table_name  => 'cred_min_apply_data_v',
  6      case_id_column_name => 'CUSTID',
  7      target_column_name => 'DEFAULTNM',
  8      settings_table_name => 'cred_model_settings');
  9      END;
 10      /
```

```
DMU17 >
DMU17 >REM Create the model using the specified settings
DMU17 > BEGIN
  2      DBMS_DATA_MINING.CREATE_MODEL(
  3      model_name      => 'cred_model',
  4      mining_function  => dbms_data_mining.classification,
  5      data_table_name  => 'cred_min_apply_data_v',
  6      case_id_column_name => 'CUSTID',
  7      target_column_name => 'DEFAULTNM',
  8      settings_table_name => 'cred_model_settings');
  9      END;
 10      /

PL/SQL procedure successfully completed.
```

```
DMU17 >REM applying the model
DMU17 >REM Find the 10 customers (display their names and IDs)
DMU17 >REM first we create the view to identify those top 10 customers
DMU17 >CREATE VIEW TOP10_Customers AS
  2      SELECT CUSTID
  3      FROM (
```

```
4      SELECT CUSTID,
5          RANK() OVER (ORDER BY PREDICTION_PROBABILITY(cred_model, 1 USING
*) DESC, CUSTID) AS RNK
6      FROM cred_min_apply_data_v
7  )
8  WHERE RNK <= 10
9      ORDER BY RNK;
```

```
DMU17 >
DMU17 >REM applying the model
DMU17 >REM Find the 10 customers (display their names and IDs)
DMU17 >REM first we create the view to identify those top 10 customers
DMU17 >CREATE VIEW TOP10_Customers AS
2  SELECT CUSTID
3  FROM (
4      SELECT CUSTID,
5          RANK() OVER (ORDER BY PREDICTION_PROBABILITY(cred_model, 1 USING *) DESC, CUSTID) AS RNK
6      FROM cred_min_apply_data_v
7  )
8  WHERE RNK <= 10
9      ORDER BY RNK;

View created.
```

```
DMU17 >REM now we display their details
DMU17 >SELECT CUSTID, ATTRB14 AS ANNUAL_INCOME
2  FROM GLOBALCREDITCARDS
3  WHERE CUSTID IN (SELECT * FROM TOP10_Customers);
```

```
DMU17 >REM now we display their details
DMU17 >SELECT CUSTID, ATTRB14 AS ANNUAL_INCOME
2  FROM GLOBALCREDITCARDS
3  WHERE CUSTID IN (SELECT * FROM TOP10_Customers);

  CUSTID ANNUAL_INCOME
-----
259319      72343
253914      42613
255241      89190
263790      29730
280421      56487
288843      42613
286421      89190
286311      89190
283344     123875
289411      56487

10 rows selected.

DMU17 >
```

```
DMU17 >REM Create the settings table
DMU17 >CREATE TABLE navbyes_model_settings (
2      setting_name VARCHAR2(30),
3      setting_value VARCHAR2(30));
```

```
DMU17 >
DMU17 >REM Create the settings table
DMU17 >CREATE TABLE navbytes_model_settings (
2      setting_name VARCHAR2(30),
3      setting_value VARCHAR2(30));

Table created.
```

```
DMU17 >REM Populate the settings table
DMU17 >BEGIN
2      INSERT INTO  navbytes_model_settings VALUES
3          (dbms_data_mining.algo_name,
4           dbms_data_mining.ALGO_NAIVE_BAYES);
5      INSERT INTO  navbytes_model_settings VALUES
6          (dbms_data_mining.prep_auto,
7           dbms_data_mining.prep_auto_on);
8      COMMIT;
9      END;
10     /
```

```
DMU17 >REM Populate the settings table
DMU17 >BEGIN
2      INSERT INTO  navbytes_model_settings VALUES
3          (dbms_data_mining.algo_name,
4           dbms_data_mining.ALGO_NAIVE_BAYES);
5      INSERT INTO  navbytes_model_settings VALUES
6          (dbms_data_mining.prep_auto,
7           dbms_data_mining.prep_auto_on);
8      COMMIT;
9      END;
10     /

PL/SQL procedure successfully completed.
```

```
DMU17 >REM Create the model using the specified settings
DMU17 >BEGIN
2      DBMS_DATA_MINING.CREATE_MODEL(
3          model_name          => 'navbytes_cred_model',
4          mining_function     => dbms_data_mining.classification,
5          data_table_name     => 'cred_min_apply_data_v',
6          case_id_column_name => 'CUSTID',
7          target_column_name  => 'DEFAULTNM',
8          settings_table_name => 'navbytes_model_settings');
9      END;
10     /
```

```
DMU17 >
DMU17 >REM Create the model using the specified settings
DMU17 >BEGIN
  2      DBMS_DATA_MINING.CREATE_MODEL(
  3          model_name          => 'navbytes_cred_model',
  4          mining_function      => dbms_data_mining.classification,
  5          data_table_name      => 'cred_min_apply_data_v',
  6          case_id_column_name  => 'CUSTID',
  7          target_column_name   => 'DEFAULTNM',
  8          settings_table_name  => 'navbytes_model_settings');
  9      END;
 10  /

PL/SQL procedure successfully completed.
```

```
DMU17 >
DMU17 >CREATE VIEW TOP10_def_Customers_nb AS
  2      SELECT CUSTID
  3      FROM (SELECT CUSTID, rank() over (order by
  4                  PREDICTION_PROBABILITY(navbytes_cred_model, 1 USING *)
  5                  DESC, CUSTID) RNK
  6              FROM cred_min_apply_data_v)
  7      WHERE RNK <= 10
  8      ORDER BY RNK;
```

```
DMU17 >
DMU17 >CREATE VIEW TOP10_def_Customers_nb AS
  2      SELECT CUSTID
  3      FROM (SELECT CUSTID, rank() over (order by
  4                  PREDICTION_PROBABILITY(navbytes_cred_model, 1 USING *)
  5                  DESC, CUSTID) RNK
  6              FROM cred_min_apply_data_v)
  7      WHERE RNK <= 10
  8      ORDER BY RNK;

View created.
```

```
DMU17 >REM now we shall display the customers names
DMU17 >SELECT CUSTID, ATTRB14 AS ANNUAL_INCOME
  2      FROM GLOBALCREDITCARDS
  3      WHERE CUSTID IN (SELECT * FROM TOP10_def_Customers_nb);
```



```
DMU17 >REM now we shall display the customers' names
DMU17 >SELECT CUSTID, ATTRB14 AS ANNUAL_INCOME
2      FROM GLOBALCREDITCARDS
3      WHERE CUSTID IN (SELECT * FROM TOP10_def_Customers_nb);
```

CUSTID	ANNUAL_INCOME
253914	42613
256565	89190
262041	29730
272265	29730
282123	72343
282592	42613
287307	72343
286482	89190
283617	42613
285762	42613

10 rows selected.

(10 marks)

4. Using suitable metrics, evaluate capabilities of the models you have developed for this task.

(8 marks)

Provide whatever PL/SQL API and SQL code you have used for this part as plain text and their outputs as screenshots. Choose a range of different evaluation metrics suitable for your data mining models.

```
DMU17 >REM testing the model
DMU17 >
DMU17 >SELECT DEFAULTNM AS actual_target_value,
2          PREDICTION(cred_model USING *) AS predicted_target_value,
3          COUNT(*) AS total_value
4      FROM cred_min_apply_data_v
5      GROUP BY DEFAULTNM, PREDICTION(cred_model USING *)
6      ORDER BY 1, 2;
```

```
DMU17 >
DMU17 >REM testing the model
DMU17 >
DMU17 >SELECT DEFAULTNM AS actual_target_value,
2          PREDICTION(cred_model USING *) AS predicted_target_value,
3          COUNT(*) AS total_value
4      FROM cred_min_apply_data_v
5      GROUP BY DEFAULTNM, PREDICTION(cred_model USING *)
6      ORDER BY 1, 2;
```

ACTUAL_TARGET_VALUE	PREDICTED_TARGET_VALUE	TOTAL_VALUE
0	0	28549
0	1	1405
1	0	6922
1	1	3124

```
DMU17 >_
```

```
DMU17 >REM calculating the model's accuracy
DMU17 > COLUMN ACCURACY FORMAT 99.99
DMU17 > SELECT (SUM(correct)/COUNT(*))*100 AS accuracy
2      FROM (SELECT DECODE(DEFAULTNM,
3          PREDICTION(cred_model USING *), 1, 0) AS correct
```



```
4          FROM cred_min_apply_data_v);

DMU17 > REM calculating the model's accuracy
DMU17 > COLUMN ACCURACY FORMAT 99.99
DMU17 > SELECT (SUM(correct)/COUNT(*))*100 AS accuracy
2          FROM (SELECT DECODE(DEFAULTNM,
3          PREDICTION(cred_model USING *), 1, 0) AS correct
4          FROM cred_min_apply_data_v);

ACCURACY
-----
79.18
```

```
DMU17 > REM testing the model
DMU17 >
DMU17 > SELECT DEFAULTNM AS actual_target_value,
2          PREDICTION(navbytes_cred_model USING *) AS
predicted_target_value,
3          COUNT(*) AS total_value
4          FROM cred_min_apply_data_v
5          GROUP BY DEFAULTNM, PREDICTION(navbytes_cred_model USING *)
6          ORDER BY 1, 2;
```

```
DMU17 >
DMU17 > REM testing the model
DMU17 >
DMU17 > SELECT DEFAULTNM AS actual_target_value,
2          PREDICTION(navbytes_cred_model USING *) AS predicted_target_value,
3          COUNT(*) AS total_value
4          FROM cred_min_apply_data_v
5          GROUP BY DEFAULTNM, PREDICTION(navbytes_cred_model USING *)
6          ORDER BY 1, 2;

ACTUAL_TARGET_VALUE  PREDICTED_TARGET_VALUE  TOTAL_VALUE
-----
0                    0                    21543
0                    1                    8411
1                    0                    2452
1                    1                    7594
DMU17 >
```

```
DMU17 > REM calculating the model accuracy
DMU17 >
DMU17 > COLUMN ACCURACY FORMAT 99.99
DMU17 > SELECT (SUM(correct)/COUNT(*))*100 AS accuracy
2          FROM (SELECT DECODE(DEFAULTNM,
3          PREDICTION(navbytes_cred_model USING *), 1,
0) AS correct
```

```
4                                FROM cred_min_apply_data_v);

DMU17 >
DMU17 >REM calculating the model accuracy
DMU17 >
DMU17 >COLUMN ACCURACY FORMAT 99.99
DMU17 >SELECT (SUM(correct)/COUNT(*)*100 AS accuracy
2      FROM (SELECT DECODE(DEFAULTNM,
3                    PREDICTION(navbytes_cred_model USING *), 1, 0) AS correct
4                    FROM cred_min_apply_data_v);

ACCURACY
-----
72.84
```

```
DMU17 >
DMU17 >REM CREATE TABLE WITH CONFUSION MATRIX DATA
DMU17 >CREATE TABLE CRED_SVM_CONFUSION_MATRIX (
2      ACTUAL_TARGET_VALUE NUMBER,
3      PREDICTED_TARGET_VALUE NUMBER,
4      TOTAL_VALUE NUMBER
5 );
```

Table created.

```
DMU17 >
DMU17 >REM INSERT DATA
DMU17 >-- Inserting values into CRED_SVM_CONFUSION_MATRIX
DMU17 >INSERT INTO CRED_SVM_CONFUSION_MATRIX VALUES (0, 0, 28549);
```

1 row created.

```
DMU17 >INSERT INTO CRED_SVM_CONFUSION_MATRIX VALUES (1, 0, 6922);
```

1 row created.

```
DMU17 >INSERT INTO CRED_SVM_CONFUSION_MATRIX VALUES (0, 1, 1405);
```

1 row created.

```
DMU17 >INSERT INTO CRED_SVM_CONFUSION_MATRIX VALUES (1, 1, 3124);
```

1 row created.

```
DMU17 >
DMU17 >
```

```
DMU17 >REM CALCULATE PRECISION, RECALL AND F1 SCORE
DMU17 >WITH ConfusionMatrix AS (
2      SELECT
3          SUM(CASE WHEN ACTUAL_TARGET_VALUE = 1 AND
PREDICTED_TARGET_VALUE = 1 THEN TOTAL_VALUE ELSE 0 END) AS TP,
4          SUM(CASE WHEN ACTUAL_TARGET_VALUE = 0 AND
PREDICTED_TARGET_VALUE = 1 THEN TOTAL_VALUE ELSE 0 END) AS FP,
5          SUM(CASE WHEN ACTUAL_TARGET_VALUE = 0 AND
PREDICTED_TARGET_VALUE = 0 THEN TOTAL_VALUE ELSE 0 END) AS TN,
6          SUM(CASE WHEN ACTUAL_TARGET_VALUE = 1 AND
PREDICTED_TARGET_VALUE = 0 THEN TOTAL_VALUE ELSE 0 END) AS FN
7      FROM CRED_SVM_CONFUSION_MATRIX
8  )
9  SELECT
10     TP, FP, TN, FN,
11     TP / (TP + FP) AS Precision,
12     TP / (TP + FN) AS Recall,
13     2 * TP / (2 * TP + FP + FN) AS F1_Score
14  FROM ConfusionMatrix;
```

```

DMU17 >
DMU17 >REM CREATE TABLE WITH CONFUSION MATRIX DATA
DMU17 >CREATE TABLE CRED_SVM_CONFUSION_MATRIX (
2     ACTUAL_TARGET_VALUE NUMBER,
3     PREDICTED_TARGET_VALUE NUMBER,
4     TOTAL_VALUE NUMBER
5 );

Table created.

DMU17 >
DMU17 >REM INSERT DATA
DMU17 >-- Inserting values into CRED_SVM_CONFUSION_MATRIX
DMU17 >INSERT INTO CRED_SVM_CONFUSION_MATRIX VALUES (0, 0, 28549);

1 row created.

DMU17 >INSERT INTO CRED_SVM_CONFUSION_MATRIX VALUES (1, 0, 6922);

1 row created.

DMU17 >INSERT INTO CRED_SVM_CONFUSION_MATRIX VALUES (0, 1, 1405);

1 row created.

DMU17 >INSERT INTO CRED_SVM_CONFUSION_MATRIX VALUES (1, 1, 3124);

1 row created.

DMU17 >
DMU17 >
DMU17 >REM CALCULATE PRECISION, RECALL AND F1 SCORE
DMU17 >WITH ConfusionMatrix AS (
2     SELECT
3         SUM(CASE WHEN ACTUAL_TARGET_VALUE = 1 AND PREDICTED_TARGET_VALUE = 1 THEN TOTAL_VALUE ELSE 0 END) AS TP,
4         SUM(CASE WHEN ACTUAL_TARGET_VALUE = 0 AND PREDICTED_TARGET_VALUE = 1 THEN TOTAL_VALUE ELSE 0 END) AS FP,
5         SUM(CASE WHEN ACTUAL_TARGET_VALUE = 0 AND PREDICTED_TARGET_VALUE = 0 THEN TOTAL_VALUE ELSE 0 END) AS TN,
6         SUM(CASE WHEN ACTUAL_TARGET_VALUE = 1 AND PREDICTED_TARGET_VALUE = 0 THEN TOTAL_VALUE ELSE 0 END) AS FN
7     FROM CRED_SVM_CONFUSION_MATRIX
8 )
9     SELECT
10        TP, FP, TN, FN,
11        TP / (TP + FP) AS Precision,
12        TP / (TP + FN) AS Recall,
13        2 * TP / (2 * TP + FP + FN) AS F1_Score
14    FROM ConfusionMatrix;

```

TP	FP	TN	FN	PRECISION	RECALL	F1_SCORE
3124	1405	28549	6922	.689776993	.31096954	.428679245

```
DMU17 >REM CREATE TABLE WITH CONFUSION MATRIX DATA
```

```

DMU17 >CREATE TABLE CRED_NB_CONFUSION_MATRIX (
2     ACTUAL_TARGET_VALUE NUMBER,
3     PREDICTED_TARGET_VALUE NUMBER,
4     TOTAL_VALUE NUMBER
5 );

```

Table created.

```
DMU17 >
```

```
DMU17 >REM INSERT DATA
```

```
DMU17 >REM Inserting values into CRED_NB_CONFUSION_MATRIX
```

```
DMU17 >INSERT INTO CRED_NB_CONFUSION_MATRIX VALUES (0, 0, 21543);
```

1 row created.

```
DMU17 >INSERT INTO CRED_NB_CONFUSION_MATRIX VALUES (1, 0, 8411);
```

1 row created.

```
DMU17 >INSERT INTO CRED_NB_CONFUSION_MATRIX VALUES (0, 1, 2452);
```

1 row created.

```
DMU17 >INSERT INTO CRED_NB_CONFUSION_MATRIX VALUES (1, 1, 7594);
```

1 row created.

```
DMU17 >
```

```
DMU17 >
```

```
DMU17 >REM CALCULATE PRECISION, RECALL AND F1 SCORE
```

```
DMU17 >WITH ConfusionMatrix AS (
```

```
2     SELECT
3         SUM(CASE WHEN ACTUAL_TARGET_VALUE = 1 AND
PREDICTED_TARGET_VALUE = 1 THEN TOTAL_VALUE ELSE 0 END) AS TP,
4         SUM(CASE WHEN ACTUAL_TARGET_VALUE = 0 AND
PREDICTED_TARGET_VALUE = 1 THEN TOTAL_VALUE ELSE 0 END) AS FP,
5         SUM(CASE WHEN ACTUAL_TARGET_VALUE = 0 AND
PREDICTED_TARGET_VALUE = 0 THEN TOTAL_VALUE ELSE 0 END) AS TN,
6         SUM(CASE WHEN ACTUAL_TARGET_VALUE = 1 AND
PREDICTED_TARGET_VALUE = 0 THEN TOTAL_VALUE ELSE 0 END) AS FN
7     FROM CRED_NB_CONFUSION_MATRIX
8 )
9 SELECT
10    TP, FP, TN, FN,
11    TP / (TP + FP) AS Precision,
12    TP / (TP + FN) AS Recall,
13    2 * TP / (2 * TP + FP + FN) AS F1_Score
14 FROM ConfusionMatrix;
```

```

DMU17 >
DMU17 >REM CREATE TABLE WITH CONFUSION MATRIX DATA
DMU17 >CREATE TABLE CRED_NB_CONFUSION_MATRIX (
2     ACTUAL_TARGET_VALUE NUMBER,
3     PREDICTED_TARGET_VALUE NUMBER,
4     TOTAL_VALUE NUMBER
5 );

Table created.

DMU17 >
DMU17 >REM INSERT DATA
DMU17 >REM Inserting values into CRED_NB_CONFUSION_MATRIX
DMU17 >INSERT INTO CRED_NB_CONFUSION_MATRIX VALUES (0, 0, 21543);

1 row created.

DMU17 >INSERT INTO CRED_NB_CONFUSION_MATRIX VALUES (1, 0, 8411);

1 row created.

DMU17 >INSERT INTO CRED_NB_CONFUSION_MATRIX VALUES (0, 1, 2452);

1 row created.

DMU17 >INSERT INTO CRED_NB_CONFUSION_MATRIX VALUES (1, 1, 7594);

1 row created.

DMU17 >
DMU17 >
DMU17 >REM CALCULATE PRECISION, RECALL AND F1 SCORE
DMU17 >WITH ConfusionMatrix AS (
2     SELECT
3         SUM(CASE WHEN ACTUAL_TARGET_VALUE = 1 AND PREDICTED_TARGET_VALUE = 1 THEN TOTAL_VALUE ELSE 0 END) AS TP,
4         SUM(CASE WHEN ACTUAL_TARGET_VALUE = 0 AND PREDICTED_TARGET_VALUE = 1 THEN TOTAL_VALUE ELSE 0 END) AS FP,
5         SUM(CASE WHEN ACTUAL_TARGET_VALUE = 0 AND PREDICTED_TARGET_VALUE = 0 THEN TOTAL_VALUE ELSE 0 END) AS TN,
6         SUM(CASE WHEN ACTUAL_TARGET_VALUE = 1 AND PREDICTED_TARGET_VALUE = 0 THEN TOTAL_VALUE ELSE 0 END) AS FN
7     FROM CRED_NB_CONFUSION_MATRIX
8 )
9     SELECT
10    TP, FP, TN, FN,
11    TP / (TP + FP) AS Precision,
12    TP / (TP + FN) AS Recall,
13    2 * TP / (2 * TP + FP + FN) AS F1_Score
14 FROM ConfusionMatrix;

      TP      FP      TN      FN  PRECISION    RECALL    F1_SCORE
-----
      7594      2452      21543      8411 .755922755 .474476726 .583010249

DMU17 >

```

5. Present and critically discuss your findings and make recommendations to the Managing Director of GLOBAL CREDIT CARDS company.

(6 marks)

In the realm of credit assessment and risk management, predictive models like Support Vector Machines (SVM) and Naive Bayes (NB) classification algorithms are indispensable tools for decision-making processes. Assessing the performance of these models provides critical insights that guide future strategies for companies like Global Credit Cards.

The evaluation of model performance underscores notable findings and directs specific recommendations for the company's advancement. SVM demonstrates a commendable accuracy of 79.18% and exhibits a balanced precision-recall trade-off with figures of 0.689, 0.311, and 0.429 respectively. On the other hand, NB yields a slightly lower accuracy at 72.84% with precision, recall, and F1 score standing at 0.756, 0.474, and 0.583. While SVM excels in overall predictive capability, NB shows potential for improvement in identifying true positive instances of defaults.

Moreover, both models successfully identify top customers based on predicted default probabilities, enabling tailored interventions for risk management. The diverse financial profiles among the top 10 predicted default customers underscore the necessity for customized credit strategies.

Moving forward, it's imperative for Global Credit Cards to refine and enhance the NB model through additional features or parameter optimization, aiming to enhance recall and overall performance. Comprehensive feature engineering should be pursued to identify pertinent predictors for both models. Continuous monitoring mechanisms must be implemented to track model performance over time and validate against real-world data. Regular validation of model assumptions and outputs is vital to detect any discrepancies or shifts in credit trends.

Furthermore, leveraging model insights to segment customers based on credit risk profiles and tailor product offerings accordingly will enhance customer engagement

and mitigate default risks. Investment in data infrastructure and analytics capabilities is crucial for advanced modeling techniques and real-time decision-making.

Promoting interdisciplinary collaboration and knowledge sharing fosters a data-driven culture within the organization, enhancing overall efficiency and effectiveness.

Compliance with regulatory requirements and ethical standards in credit assessment and lending practices is paramount. Transparency, equity, customer privacy, and data security must be prioritized in credit evaluation processes to maintain trust and integrity.

In conclusion, the deployment of predictive models presents significant opportunities for Global Credit Cards to optimize credit risk management practices and enhance customer-centricity. By adhering to the recommendations for model refinement, customer engagement, data-driven decision-making, and compliance, the company can effectively navigate dynamic credit landscapes and sustain long-term success in the financial services industry.

Part 3 (15 marks)

Critically evaluate the SH data warehouse and the GLOBAL CREDIT CARDS company's GlobalCreditCards dataset in relation to the theory and best practices of data quality and standards.

The report should be concise and comprehensive and in the region of 900-1000 words. You should use Harvard style of citation and referencing by following the guidelines in Pears and Shields (2008).

Answer Part 3: 15 Marks [10 for the quality of your report addressing the above points, 3 for the quality of referencing and citation and adhering to the Harvard style, 2 for presentation of the report]

Introduction:

Effective management of data quality is crucial to ensuring the reliability and usability of databases. This report evaluates the SH data warehouse and the GlobalCreditCards dataset, focusing on key data quality dimensions such as accuracy, completeness, consistency, and uniqueness. The evaluation is guided by established data quality frameworks and industry best practices.

Accuracy:

Accuracy in data quality refers to the degree to which data precisely represents the actual entities or occurrences it aims to describe. Accurate data is free from errors, inconsistencies, or discrepancies that could lead to misinformation. An example in the SH data warehouse is the presence of different column names like 'CUST_ID' and 'CUSTOMER_ID,' which can potentially introduce errors in queries and impact accuracy (Redman, 1996). Addressing these discrepancies and ensuring consistent column names is recommended. Conversely, the GlobalCreditCards dataset maintains accuracy through a comprehensive data dictionary, providing precise and standardized definitions for each attribute (Wang & Strong, 1996).

Completeness:

Completeness in data quality reflects the extent to which a dataset includes all necessary data items. A complete dataset contains required information without any null or missing values. Concerns about completeness arise in the SH data warehouse, notably with 10,000 rows in the PRODUCTS database containing at least one NULL value. In contrast, the GlobalCreditCards dataset demonstrates

completeness through a well-defined data dictionary outlining expected features and values.

Consistency:

Consistency pertains to the uniformity and coherence of data within a dataset or across datasets. Consistent data minimizes variations by ensuring that similar information is represented in a standardized manner. The SH data warehouse encounters consistency issues, including duplicate or unnecessary columns in the TIMES table and duplicate entries in critical tables (SALES, COSTS, and PRODUCTS). The GlobalCreditCards dataset maintains consistency through a well-organized data dictionary, reducing the likelihood of inconsistencies.

Uniqueness:

Uniqueness in data quality ensures that every record or data element in a dataset is unique and not duplicated. Duplicate entries can lead to misrepresentation of underlying entities and introduce errors. Uniqueness concerns arise in the SH data warehouse, evidenced by duplicate entries in the SALES, COSTS, and PRODUCTS tables. The GlobalCreditCards dataset is not explicitly mentioned to have uniqueness issues, suggesting effective management in this regard.

By adhering to these principles, databases can enhance their overall data quality and contribute to more reliable and effective data-driven decision-making processes.

References & Bibliography

Northumbria (2023) Academic Regulations for Taught Awards (ARTA) – 2022/23. Available at: <https://www.northumbria.ac.uk/about-us/university-services/student-library-and-academic-services/quality-and-teaching-excellence/assessment/assessment-regulations-and-policies/> (Accessed: 30 September 2023).

Pears, R. and Shields, G. (2008) *Cite them right: the essential referencing guide*. Newcastle upon Tyne: Pear Tree Books. Available at: <https://www.citethemrightonline.com/> (Accessed: 16 October 2023).

Potineni, P. (2021) Oracle Database Data Warehousing Guide, 19c. Part Number E96243-05. Available at: <https://docs.oracle.com/en/database/oracle/oracle-database/19/dwhsg/> (Accessed: 25 October 2023).

Surampudi, S. (2017a) Data Mining User's Guide, 12c Release 1 (12.1). Part Number E53115-05. Available at: <https://docs.oracle.com/database/121/DMPRG/toc.htm> (Accessed: 16 October 2023).

Surampudi, S. (2017b) Oracle Data Mining Concepts, 12c Release 1 (12.1). Part Number E17692-19. Available at: <https://docs.oracle.com/database/121/DMCON/toc.htm> (Accessed: 16 October 2023).

References and Bibliography:

peter, O. (1997) *Improved query performance with variant indexes - semantic scholar, Improved-query-performance-with-variant-indexes*. Available at: <https://www.semanticscholar.org/paper/Improved-query-performance-with-variant-indexes-O'NeilPatrick-QuassDallan/d07cdad38ca019b13ee67d233634a35b2b291636> (Accessed: 28 January 2024).

corporation, oracle (2023) *Oracle Database 19C - Data Warehousing, Oracle Help Center*. Available at: <https://docs.cloud.oracle.com/en/database/oracle/oracle-database/19/data-warehousing.html> (Accessed: 28 January 2024).

Ramkrishnan, R. and Gehrke, J. (2003) *Database Management Systems, Database Management Systems (Third Edition)*. Available at: <https://pages.cs.wisc.edu/~dbbook/> (Accessed: 28 January 2024).

corporation, oracle (2023a) *Oracle Database 19C - Data Warehousing, Oracle Help Center*. Available at: <https://docs.cloud.oracle.com/en/database/oracle/oracle-database/19/data-warehousing.html> (Accessed: 28 January 2024).

Kimball, R. and Ross, M. (2002) *The Data Warehouse Toolkit, 3rd Edition, Kimball Group*. Available at: <https://www.kimballgroup.com/data-warehouse-business-intelligence-resources/books/data-warehouse-dw-toolkit/> (Accessed: 28 January 2024).

Dam, S. (2009) *SQL Server 2008 query performance tuning distilled, Google Books*. Available at: https://books.google.com/books/about/SQL_Server_2008_Query_Performance_Tuning.html?id=GZJQUi2Yx-YC (Accessed: 28 January 2024).

Wang, R.Y., Technology, M.I. of and Metrics, O.M.A. (1998) *A product perspective on total data quality management, Communications of the ACM*. Available at: <https://dl.acm.org/doi/abs/10.1145/269012.269022> (Accessed: 28 January 2024).