

Master of Computer and Information Sciences

Assessment 2 – Version 02 – DW Project



Building and Analysing a DW for Countdown Stores in NZ

Paper: Data Warehousing and Big Data

Paper Code: COMP810

Semester-2, 2015

Weight in your final grade: 80%

Note: To pass the paper, student needs to attempt both assessments and obtain a C- grade overall.

1. Assessment task

Student has to design, implement and analyse a Data Warehouse (DW) for Countdown, one of the biggest super market chains in NZ.

2. Project overview

Countdown is a one of the biggest superstores chains in NZ. The stores locate all over the country. Countdown has thousands of customers and therefore it is important for the organisation to analyse the shopping behaviour of their customers. Based on that the organisation can optimise their selling techniques e.g. giving of promotions on different products.

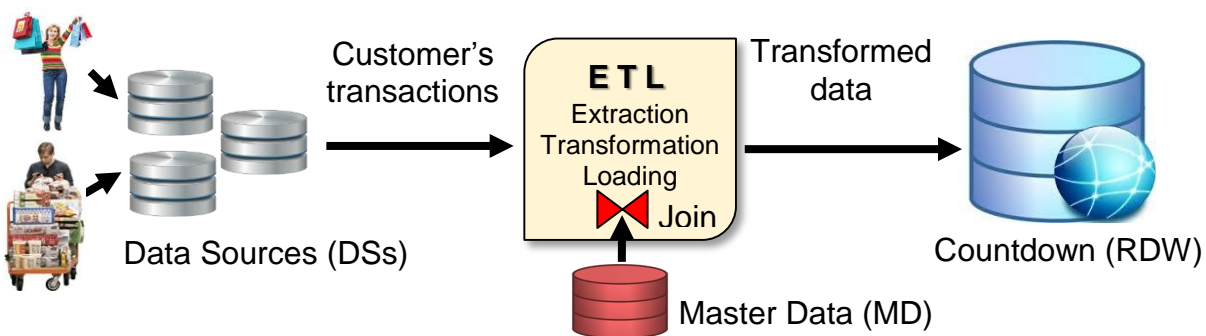


Figure 1: An overview of data integration in countdown scenario

Now, to make this analysis of shopping behaviour practical there is a need of building a DW and customers' transactions from Data Sources (DS) are required to reflect into DW on daily basis. This process of reflecting the customer data into DW is called Data Integration (DI) as shown in Figure 1. To implement DI we usually need ETL (Extraction, Transformation, and Loading) tools. Since the data generated by customers is not in the format required by DW therefore, it needs to process in the transformation layer of ETL. For example enriching of some information from Master Data (MD) as shown in Figure 2.

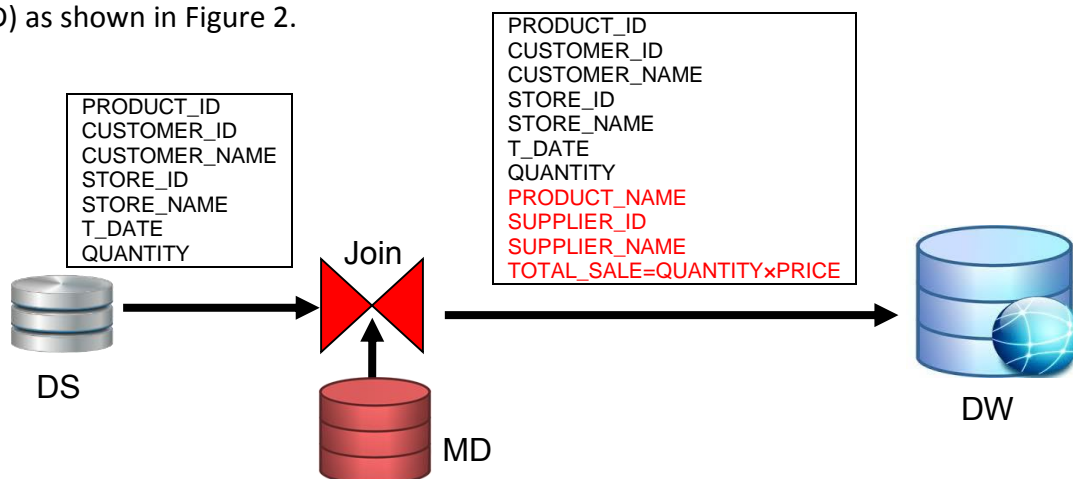


Figure 2: Enrichment example

To implement this enrichment feature in the transformation phase of ETL we need a join operator typically called Semi-Stream Join (SSJ). There are a number of algorithms available to implement this join operation however, the simplest one is Index Nested Loop Join (INLJ) which is explained in next section and you will implement it in this project.

3. Index Nested Loop Join (INLJ)

The INLJ is a traditional join operator to implement the join operation between DS and MD. In INLJ, DS is scanned tuple by tuple and based on that the disk-based MD is accessed using a cluster-based index on the join attribute. A graphical overview of an INLJ is shown in Figure 3 where DS tuple s_i is joined with MD tuple r_j .

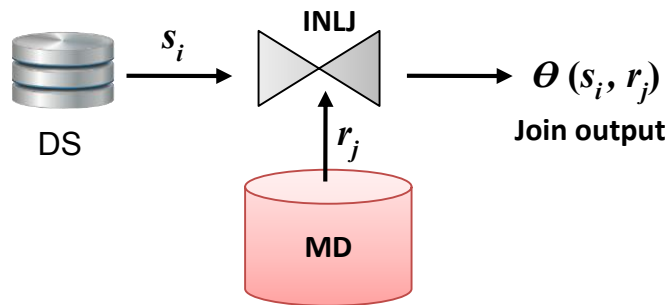


Figure 3: Execution architecture of INLJ

The crux of INLJ is that with every loop step the algorithm reads a chunk of DS tuples and joins them one-by-one with the relevant tuples in MD. To access the relevant tuple from MD the algorithm uses index on the join attribute in MD. For example, for the scenario presented in Figure 2 the join attribute is `PRODUCT_ID` which should be considered as a Primary Key (PK) in MD while a Foreign Key (FK) in DS and therefore, the join would be between PK and FK.

4. Star-schema

The star schema (which you will use in this project) is a data modelling technique that is used to map multidimensional decision support data into a relational database. Star-schema yields an easily implemented model for multidimensional data analysis while still preserving the relational structures on which the operational database is built.

The star schema represents aggregated data for specific business activities. Using the schema, one can create multiple aggregated data sources that will represent different aspects of business operations. For example, the aggregation may involve total sales by selected time periods, by products, by stores, and so on. Aggregated totals can be total product sold, total sales values by products, etc. The basic star schema has three main components: *facts*, *dimensions*, *attributes*. Usually in case of star-schema for sales the dimension tables are: *product*, *date*, *store*, and *supplier* while the fact table is *sales*. However, to determine the right attributes you will consider Figure 2.

5. Data specifications

The assessment provides a scripts file named “Transaction_and_MasterData_Generator.sql”. By executing the script it will create two tables in your account. One is `TRANSACTIONS` table with 10,000 records populated in it. This data will be generated randomly based on 100 products, 50 customers, 10 stores, and one year time period as a date - from 01-Jan-14 to 31-Dec-14. The values for the quantity attribute will be random between 1 and 10. The other is `MASTERDATA` table with 100 records in it. The structure of both tables with their attributes name and data types

is given below in Figure 4. The attributes TRANSACTION_ID and PRODUCT_ID are primary keys in TRANSACTIONS and MASTERDATA tables respectively.

TRANSACTIONS								
Attributes	TRANSACTION_ID	PRODUCT_ID	CUSTOMER_ID	CUSTOMER_NAME	STORE_ID	STORE_NAME	T_DATE	QUANTITY
Data type and size	VARCHAR2(8)	VARCHAR2(6)	VARCHAR2(4)	VARCHAR2(30)	VARCHAR2(3)	VARCHAR2(20)	DATE	NUMBER(3,0)

MASTERDATA					
Attributes	PRODUCT_ID	PRODUCT_NAME	SUPPLIER_ID	SUPPLIER_NAME	PRICE
Data type and size	VARCHAR2(6)	VARCHAR2(30)	VARCHAR2(5)	VARCHAR2(30)	NUMBER(5,2) DEFAULT 0.0

Figure 4: Structures for TRANSACTIONS and MASTERDATA tables

6. Implementation of INLJ

To implement INLJ algorithm you will implement the following steps.

1. Read 50 tuples from TRANSACTIONS table as an input data into a cursor. The cursor is a user defined data type in PLSQL which works as a list and uses to store multiple records in memory.
2. Read the cursor tuple by tuple and for each tuple retrieve the relevant tuple from MASTERDATA table using PRODUCT_ID as an index and add the required attributes (mentioned in Figure 2) into the transaction tuple.
3. The transaction tuple with new attributes will be loaded into DW. Before loading the tuple into DW you will check whether the dimensions tables already contains this information. If yes then only update the fact table otherwise update both dimensions and fact tables.
4. Repeat steps 1 to 3 until you load all the data from TRANSACTIONS table to DW.

7. DW analysis

Once the entire data has been loaded into DW, apply the following analysis to your DW using OLAP queries.

1. Which product generated maximum sales in Dec, 2014?
2. Which store produced highest sales in the whole year?
3. Determine the supplier name for the most popular product based on sales.
4. Presents the quarterly sales analysis for all stores using drill down query concepts.
5. Create a view with name "STOREANALYSIS" that present the product-wise sales analysis for each store.

8. Tasks break-up

Following is a list of tasks that you need to complete in this project.

1. Identifying of appropriate dimension tables, fact table, and their attributes for the sales scenario presented in Figure 2. Based on that creating of star-schema for DW with appropriate primary and foreign keys. To keep the attribute name and their data types consistent in DW, consult the structure of tables TRANSACTIONS and MASTERDATA provided in Figure 4.
2. Implementing of the INLJ algorithm and loading of transactional data into DW after joining it with MD.
3. Applying of different analysis (described in Section 7) on DW using slicing, dicing, drill down, and materialising view concepts.
4. Writing of project report that should include project overview, INLJ algorithm, schema for DW, your OLAP queries with outputs and what did you learn from the project?

9. What to submit

Each student has to submit the following files:

1. *createDW* –SQL script file to create star-schema for DW
Note: your scripts should drop the table(s) if they already exist in the database.
2. *INLJ* – PLSQL file that implements the INLJ algorithm
3. *queriesDW* – SQL script file containing all of your OLAP queries
4. *projectReport* – a doc file containing all contents described in point 4 under the task break-up section
5. *readMe* – a text file describing the step-by-step instructions to operate your project

Note: all above files need to submit in a zipped folder named by your family name, student ID, and assessment version e.g. John-12345v2.

10. When to submit

Due date: **Friday, 30th Oct 2015**

Late penalty: maximum late submissions time is 24 hours after the due date. In this case 5% late penalty will be applied.

11. Who to submit

The project should be submitted through autonline.

NOTE: Every student has to complete the project individually. Each student's project source and report materials should be unique and done by his/her own. All assessments will be assessed through turnitin system and in case of finding of any duplication or identical material the AUT cheating policy will be applied.

-----E N D-----

Marking guide

Project Component	Marks
<i>createRDW</i> –SQL script file to create star-schema for RDW	/15
The script should create all dimensions' and fact tables table in RDW and if any table with same name exists already the script should drop that. It should also apply all primary and foreign keys on the right attributes.	
Implementing of INLJ	/30
INLJ procedure should implement all three phases of ETL – it should extract records from TRANSACTIONS table, transform these with MD and then load these records to DW successfully.	
<i>queriesDW</i> – SQL script file containing of all your OLAP queries	/35
The file should include OLAP queries for all tasks presented in Section 7.	
<i>projectReport</i> – a doc file containing all contents described in point 4 under the task break-up section.	/15
Report must contains project overview, INLJ algorithm, schema for DW, your OLAP queries with outputs and what did you learn from the project?	
<i>readMe</i> – a text file describing the step-by-step instructions to operate your project	/5
readMe file should contain a step-by-step guide to operate the project.	
Late submission penalty	-/5
TOTAL MARKS	/100