Class: BE 3

Batch: P3

Roll no. 41310

Name: Prem Vinod Bansod

**Date: 19-8-2020**

# Assignment 1

## Problem Definition:

For an organization of your choice, choose a set of business processes. Design star / snow flake schemas for analysing these processes. Create a fact constellation schema by combining them. Extract data from different data sources, apply suitable transformations and load into destination tables using an ETL tool.

For Example: Business Origination: Sales, Order, Marketing Process.

## Learning Objective:

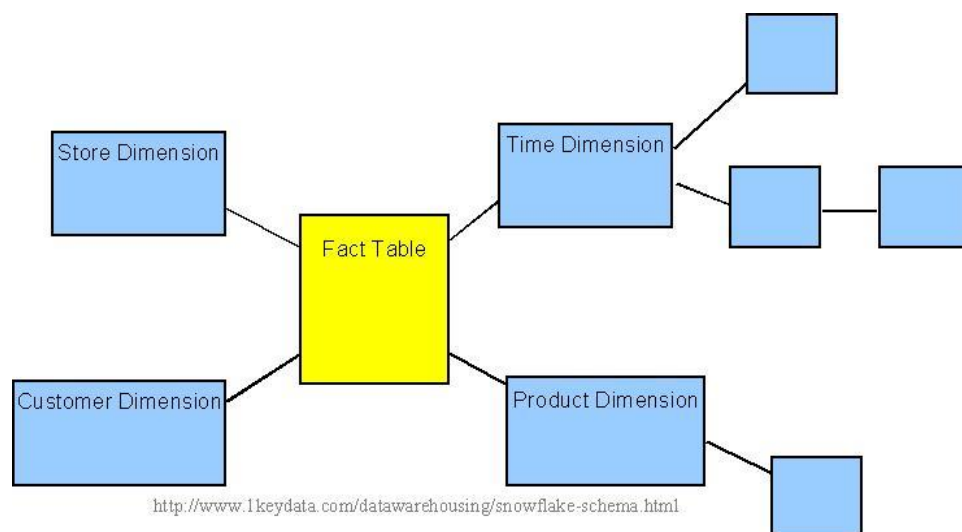- Implementation of the problem statement using ETL Tool.
- Star / snow flake schemas for analysing processes.

## Theory:

**Star / snow flake schemas:**

The star schema is the simplest type of Data Warehouse schema. It is known as star schema as its structure resembles a star. A Snowflake Schema is an extension of a Star Schema, and it adds additional

dimensions. It is called snowflake because its diagram resembles a Snowflake. Star schema and the snowflake schema are ways to organize data marts or entire data warehouses using relational databases. The third differentiator in this Star schema vs Snowflake schema face-off is the performance of these models. The Snowflake model has more joins between the dimension table and the fact table, so the performance is slower. The Star model, on the other hand, has fewer joins between dimension tables and the facts table.



http://www.1keydata.com/datawarehousing/snowflake-schema.html

**Characteristics of Star Schema:**

1. The dimension table should contain the set of attributes.
2. The dimension table is joined to the fact table using a foreign key.
3. The dimension table are not joined to each other.
4. The schema is widely supported by BI Tools.

**Snowflake Schema:**

Snowflake Schema is also the type of multidimensional model which is used for data warehouse. In snowflake schema, the fact tables, dimension tables as well as sub dimension tables are contained. This schema forms a snowflake with fact tables, dimension tables as well as sub- dimension tables.

**Characteristics of Snowflake Schema:**

1. The main benefit of the snowflake schema it uses smaller disk space.
2. Easier to implement a dimension is added to the Schema.
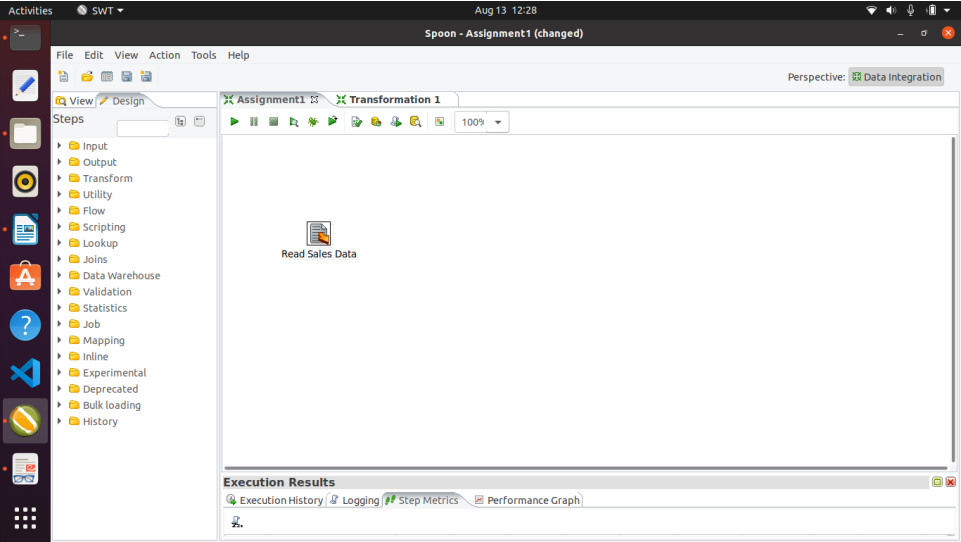3. Due to multiple tables query performance is reduced.

**What is ETL?**

ETL is an abbreviation of Extract, Transform and Load. In this process, an ETL tool extracts the data from different RDBMS source systems then transforms the data like applying calculations, concatenations, etc. and then load the data into the Data Warehouse system. In ETL data is flows from the source to the target. In ETL process transformation engine takes care of any data changes.

**List of open sources ETL Tools:**

1. CloverETL
2. Jedox
3. Pentaho
4. Talend

## Steps:

**Scan results**

Result:

Result after scanning 5647 lines.
------------------------------------------------
Field nr. 1 :
  Field name          : ORDERNUMBER
  Field type          : Integer

Field nr. 2 :
  Field name          : QUANTITYORDERED
  Field type          : Integer

Field nr. 3 :
  Field name          : PRICEEACH
  Field type          : Number
  Estimated length    : 15
  Estimated precision  : -
  Number format        : #.#
    WARNING: More then 1 number format seems to match all sampled records:
    Number format      : #.#
      Trim Type        : 0
      Minimum value    : 26.88
      Maximum value    : 100.0
      Example          : #.#, number [26.88] gives 26.88
    Number format      : #.0
      Trim Type        : 0
      Minimum value    : 26.0
      Maximum value    : 100.0
      Example          : #.#, number [26.0] gives 26.0
    Number format      : #.0
      Trim Type        : 0
      Minimum value    : 26.88
      Maximum value    : 100.0
      Example          : #.0, number [26.88] gives 26.88
    Number format      : #.0
      Trim Type        : 0
      Minimum value    : 26.0
      Maximum value    : 100.0
      Example          : #.0, number [26.0] gives 26.0
    Number format      : #.#
      Trim Type        : 3
      Minimum value    : 26.88
      Maximum value    : 100.0
      Example          : #.#, number [26.88] gives 26.88
    Number format      : #.0
      Trim Type        : 3
      Minimum value    : 26.0
      Maximum value    : 100.0
      Example          : #.#, number [26.0] gives 26.0

Close

---

File | Content | Error Handling | Filters | Fields | Additi...

| # | Name | Type | | | mal | Group | Null if | Default | Trim type |
|---|------|------|---|---|-----|-------|---------|---------|-----------|
| 4 | ORDERLINENUMBER | Integer | | | | , | | - | none |
| 5 | SALES | Number | | | | , | | - | none |
| 6 | ORDERDATE | Date | | | | | | - | none |
| 7 | STATUS | String | | | | | | - | none |
| 8 | QTR_ID | Integer | | | | | | - | none |
| 9 | MONTH_ID | Integer | | | | | | - | none |
| 10 | YEAR_ID | Integer | | | | | | - | none |
| 11 | PRODUCTLINE | String | | | | | | - | none |
| 12 | MSRP | Integer | | | | | | - | none |
| 13 | PRODUCTCODE | String | | | | | | - | none |
| 14 | CUSTOMERNAME | String | | | | | | - | none |
| 15 | PHONE | String | | | | | | - | none |
| 16 | ADDRESSLINE1 | String | | | | | | - | none |
| 17 | ADDRESSLINE2 | String | | | | | | - | none |
| 18 | CITY | String | | | | | | - | none |
| 19 | STATE | String | | | | | | - | none |
| 20 | POSTALCODE | String | | | | | | - | none |
| 21 | COUNTRY | String | | | | | | - | none |
| 22 | TERRITORY | String | | | | | | - | none |
| 23 | CONTACTLASTNAME | String | | | | | | - | none |
| 24 | CONTACTFIRSTNAME | String | | | | | | - | none |
| 25 | ORDERNUMBER | Integer | | | | , | | | none |

---

**Spoon - Assignment1 (changed)**

File  Edit  View  Action  Tools  Help

Perspective: Data Integration

View | Design

Steps

- Input
- Output
- Transform
- Utility
- Flow
- Scripting
- Lookup
- Joins
- Data Warehouse
- Validation
- Statistics
- Job
- Mapping
- Inline
- Experimental
- Deprecated
- Bulk loading
- History

Assignment1 | Transformation 1

100%

Read Sales Data → Filter Missing zips

**Execution Results**

Execution History | Logging | Step Metrics | Performance Graph

**Content of first file**

100 lines:

```
CITY,STATE,POSTALCODE
ABBEVILLE,AL,36310
ABBEVILLE,LA,70510
ABBEVILLE,MS,38601
ABBOT,ME,4406
ABBOTT,TX,76621
ABBYVILLE,KS,67510
ABERCROMBIE,ND,58001
ABERDEEN,KY,42201
ABERDEEN,MS,39730
ABERDEEN,OH,45101
ABERDEEN,SD,57402
ABERDEEN PROVING GROUND,MD,21005
ABERNATHY,TX,79311
ABILENE,KS,67410
ABILENE,TX,79602
ABILENE,TX,79604
ABILENE,TX,79606
ABILENE,TX,79697
ABILENE,TX,79699
ABINGDON,MD,21009
ABINGDON,VA,24211
ABINGTON,CT,6230
ABINGTON,PA,19001
ABITA SPRINGS,LA,70420
```

Close

**Examine preview data**

Rows of step: Read Postal Codes (1000 rows)

| # | CITY | STATE | POSTALCODE |
|---|------|-------|------------|
| 1 | ABBEVILLE | AL | 36310 |
| 2 | ABBEVILLE | LA | 70510 |
| 3 | ABBEVILLE | MS | 38601 |
| 4 | ABBOT | ME | 4406 |
| 5 | ABBOTT | TX | 76621 |
| 6 | ABBYVILLE | KS | 67510 |
| 7 | ABERCROMBIE | ND | 58001 |
| 8 | ABERDEEN | KY | 42201 |
| 9 | ABERDEEN | MS | 39730 |
| 10 | ABERDEEN | OH | 45101 |
| 11 | ABERDEEN | SD | 57402 |
| 12 | ABERDEEN PROVING GROUND | MD | 21005 |
| 13 | ABERNATHY | TX | 79311 |
| 14 | ABILENE | KS | 67410 |
| 15 | ABILENE | TX | 79602 |
| 16 | ABILENE | TX | 79604 |
| 17 | ABILENE | TX | 79606 |
| 18 | ABILENE | TX | 79697 |
| 19 | ABILENE | TX | 79699 |
| 20 | ABINGDON | MD | 21009 |
| 21 | ABINGDON | VA | 24211 |
| 22 | ABINGTON | CT | 6230 |
| 23 | ABINGTON | PA | 19001 |
| 24 | ABITA SPRINGS | LA | 70420 |
| 25 | ABSARAKA | ND | 58002 |
| 26 | ABSECON | NJ | 8201 |

Close    Show Log

**Spoon - Assignment1 (changed)**

File   Edit   View   Action   Tools   Help

Perspective: Data Integration

View / Design     Assignment1   Transformation 1

Steps

- Input
- Output
- Transform
- Utility
- Flow
- Scripting
- Lookup
- Joins
- Data Warehouse
- Validation
- Statistics
- Job
- Mapping
- Inline
- Experimental
- Deprecated
- Bulk loading
- History

100%

Read Sales Data    Filter Missing zips    Write to database

Read Postal Codes    Lookup missing zips    Prepare field layout

**Execution Results**

Execution History   Logging   Step Metrics   Performance Graph

## Conclusion:

Thus, Successfully learned how to extract data from different data sources using pentaho and apply suitable transformations and load into destination tables using an ETL tool.