**Assignment A1 ( Data Mining & Warehousing )**

**Problem Definition**:
For an organization of your choice, choose a set of business processes. Design star / snowflake schemas for analyzing 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.

**Requirements**:
●Windows 10 ·
●Open source ETL tool for Linux/Windows - Pentaho

**Learning Objectives**:

●To understand the concept of Data cube and multidimensional model for data warehouse.
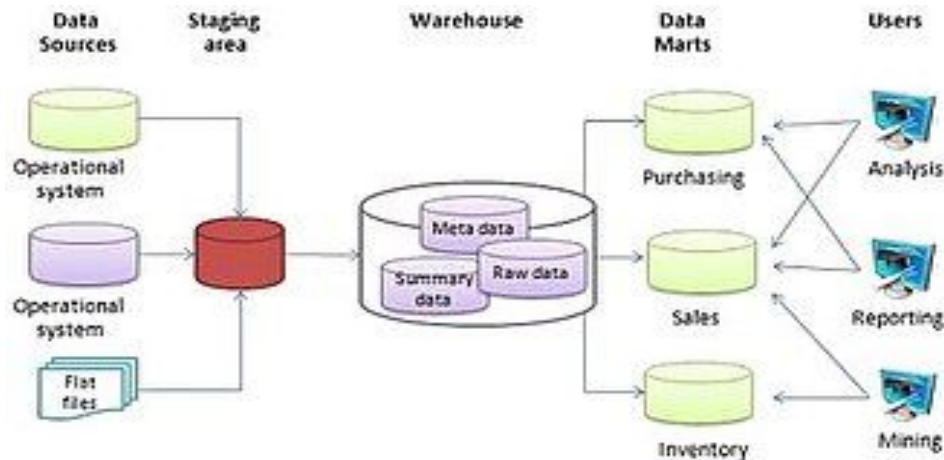●To understand different preprocessing techniques.
●To study ETL tool Pentaho.

**Learning Outcomes**:

●Understood the concept of data cube and multidimensional model for data warehouse.
●Applied preprocessing technique in order to replace missing data in the dataset.
●Used pentaho ETL tool to extract data from sales data set, apply suitable transformations and load into destination tables.

**Theory**:
- **Data Warehouse**
  A data warehouse is a central repository of information that can be analyzed to make more informed decisions. Data flows into a data warehouse from transactional systems, relational databases, and other sources, typically on a regular cadence. Business analysts, data engineers, data scientists, and decision makers access the data through business intelligence (BI) tools, SQL clients, and other analytics applications.A data warehouse is usually modeled by a multidimensional data structure, called a data cube, in which each dimension corresponds to an attribute or a set of attributes in the schema, and each cell stores the value of some aggregate measure such as count,or sum(sales amount).  A data cube provides a multidimensional view of data and allows the pre-computation and fast access of summarized data. A multidimensional model views data in the form of a data-cube. A data cube enables data to be modeled and viewed in multiple dimensions. It is defined by dimensions and facts.

**Overview**

- **Data Warehousing**
  Data warehousing is the process of constructing and using a data warehouse. A data warehouse is constructed by integrating data from multiple heterogeneous sources that support analytical reporting, structured and/or ad hoc queries, and decision making. Data warehousing involves data cleaning, data integration, and data consolidations.

  Data Warehousing allows for the following:
  - Tuning Strategies
  - Customer analysis
  - Product analysis

  In order to derive value from your data, you need not only have it in one place and using a single, canonical access language and interface, but you must also have a means to manage metadata, handle governance issues, and scale as your data grows. These are among the many challenges that data warehousing tools solve.

  Functions of Data Warehousing tools:
  - Data Extraction − Involves gathering data from multiple heterogeneous sources.

  - Data Cleaning − Involves finding and correcting the errors in data.

  - Data Transformation − Involves converting the data from legacy format to warehouse format.

  - Data Loading − Involves sorting, summarizing, consolidating, checking integrity, and building indices and partitions.

- Refreshing − Involves updating from data sources to warehouse

- **ETL Data Warehousing Tool**
  ETL stands for "Extract, Transform, and Load" and consists of the tools and processes used for pulling data from one store, transforming it for placement, and finally, loading it into another (often aggregate) store. Just as with data warehouses, ETL tools have progressed over time from self-administered to cloud-native offerings.

- **Star Schema**
  Star schema is a mature modeling approach widely adopted by relational data warehouses. It requires modelers to classify their model tables as either dimension or fact. Dimension tables describe business entities—the things you model. Entities can include products, people, places, and concepts including time itself. The most consistent table you'll find in a star schema is a date dimension table. A dimension table contains a key column (or columns) that acts as a unique identifier, and descriptive columns. Fact tables store observations or events, and can be sales orders, stock balances, exchange rates, temperatures, etc. A fact table contains dimension key columns that relate to dimension tables, and numeric measure columns. The dimension key columns determine the dimensionality of a fact table, while the dimension key values determine the granularity of a fact table.

- **Snowflake Schema**
  It is a logical arrangement of tables in a multidimensional database such that the ER diagram resembles a snowflake shape. A Snowflake Schema is an extension of a Star Schema, and it adds additional dimensions. The dimension tables are normalized which splits data into additional tables. The main benefit of the snowflake schema is that it uses smaller disk space. Easier to implement a dimension is added to the Schema
  Due to multiple tables query performance is reduced. The primary challenge that you will face while using the Snowflake Schema is that you need to perform more maintenance efforts because of the more lookup tables.

| Star Schema | Snowflake Schema |
|---|---|
| Hierarchies for the dimensions are stored in the dimensional table. | Hierarchies are divided into separate tables |
| It contains a fact table surrounded by dimension tables. | One fact table surrounded by dimension table which are in turn surrounded by dimension table |
| In a star schema, only single join creates the relationship between the fact table and any dimension tables | A snowflake schema requires many joins to fetch the data. |
| Simple DB Design. | Very Complex DB Design. |
| Denormalized Data structure and query also run faster. | Normalized Data Structure. |
| High level of Data redundancy | Very low-level data redundancy |

**OPERATIONS PERFORMED:**
1. Retrieve sales data from sample sales data csv file (Extraction).
2. Filter records with missing postal codes (POSTALCODE = null).
3. Load the data into MySQL database.
4. Perform lookup operation for missing zip codes.
5. Resolve missing zip codes.
6. Write the resolved data into the database.

Spoon - [PRODUCTION-1] Getting Transformation

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

Getting Transformation

100%

Read Sales Data    Filter Missing Zips    Table output

Select values

Read Postal Codes    Lookup Missing Zips

**Execution Results**

Execution History | Logging | Step Metrics | Performance Graph

| # | Stepname | Copynr | Read | Written | Input | Output | Updated | Rejected | Errors | Active | Time | Speed (r/s) | input/output |
|---|----------|--------|------|---------|-------|--------|---------|----------|--------|--------|------|-------------|--------------|
| 1 | Read Sales Data | 0 | 0 | 2823 | 2824 | 0 | 1 | 0 | 0 | Finished | 0.3s | 8455.0 | - |
| 2 | Filter Missing Zips | 0 | 2823 | 2823 | 0 | 0 | 0 | 0 | 0 | Finished | 0.3s | 8327.4 | - |
| 3 | Table output | 0 | 2823 | 2823 | 0 | 2823 | 0 | 0 | 0 | Finished | 8.5s | 331.1 | - |
| 4 | Read Postal Codes | 0 | 0 | 21379 | 21380 | 0 | 1 | 0 | 0 | Finished | 0.3s | 76357.1 | - |
| 5 | Lookup Missing Zips | 0 | 21455 | 76 | 0 | 0 | 0 | 0 | 0 | Finished | 0.4s | 48105.3 | - |
| 6 | Select values | 0 | 76 | 76 | 0 | 0 | 0 | 0 | 0 | Finished | 0.5s | 167.7 | - |

Type here to search    ENG  US  15:07  13-08-2020

---

MySQL Workbench

Local instance MySQL80

File  Edit  View  Query  Database  Server  Tools  Scripting  Help

Navigator

SCHEMAS

Filter objects

- r_step_type
- r_trans_attribute
- r_trans_cluster
- r_trans_hop
- r_trans_lock
- r_trans_note
- r_trans_partition_sch
- r_trans_slave
- r_trans_step_conditi
- r_transformation
- r_user
- r_value
- r_version
- sales_data
- Views
- Stored Procedures
- Functions
- sakila
- sys
- world

Administration  Schemas

Information

No object selected

sales_data

Limit to 1000 rows

```
1 •  SELECT * FROM pentaho_repo.sales_data where POSTALCODE = null
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content:

| ORDERNUMBER | QUANTITYORDERED | PRICEEACH | ORDERLINENUMBER | SALES | ORDERDATE | STATUS | QTR_ID | MONTH_ID | YEAR_ID | PRODUCTLINE | MSRP | PRODUCTCODE | CUSTOMERNAME | PHONE | ADDRESSLINE1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

sales_data3

Read Only

Output

Action Output

| # | Time | Action | Message | Duration / Fetch |
|---|------|--------|---------|------------------|
| 1 | 14:37:23 | SELECT * FROM pentaho_repo.sales_data LIMIT 0, 1000 | 1000 row(s) returned | 0.063 sec / 0.078 sec |
| 2 | 14:37:58 | SELECT * FROM pentaho_repo.sales_data where POSTALCODE = null LIMIT 0, 1000 | 0 row(s) returned | 0.109 sec / 0.000 sec |
| 3 | 15:07:30 | SELECT * FROM pentaho_repo.sales_data where POSTALCODE = null LIMIT 0, 1000 | 0 row(s) returned | 0.031 sec / 0.000 sec |

Object Info  Session

Type here to search    ENG  US  15:07  13-08-2020

**Conclusion**:

Used pentaho ETL tool to extract data from sales data set, applied preprocessing technique in order to replace missing zip codes and loaded the data into destination tables