| |
|---|
| **Vishwakarma Institute Of Information Technology, Pune.** |

Department Computer Engineering (Software Engineering)

(Data Warehouse and Data Analytics Lab)

| |
|---|
| **EXPERIMENT NO – 01** |

| |
|---|
| **Aim :** Choose a set of business processes like Sales, Customer Services, Accounting, Production, Marketing processes etc. for any organization and design star, snowflake and fact constellation schema. Also using ETL tools ,extract data from various sources and perform transform and load operations on data. (Use Power BI/RapidMiner). |

**Objective:** To design dimensional models (star, snowflake, and fact constellation schemas) for a set of business processes and implement an ETL process to populate the data warehouse using a chosen ETL tool (Power BI or RapidMiner).

**Theory :**

**Scenario:** Choose an organization (real or hypothetical) and select at least three business processes within that organization. Examples:

- Retail: Sales, Customer Service, Marketing
- Manufacturing: Production, Inventory, Sales
- Healthcare: Patient Care, Billing, Admissions
- Education: Student Enrollment, Course Management, Alumni Relations

**Tasks:**

1. **Business Process Definition:**
   o Clearly define the chosen business processes.
   o Describe the key metrics and dimensions associated with each process.
   o Identify the data sources for each process (e.g., CRM, ERP, databases, files).
2. **Dimensional Modeling:**
   o **Star Schema:** Design a star schema for *one* of the chosen business processes. Include a fact table and dimension tables. Clearly define the grain of the fact table.
   o **Snowflake Schema:** Design a snowflake schema for the *same* business process as the star schema. Normalize at least one dimension table. Explain the benefits of the snowflake schema in this context.
   o **Fact Constellation Schema:** Design a fact constellation schema that integrates *all* chosen business processes. Identify shared dimension tables. Explain how this schema supports cross-process analysis.
3. **ETL Process Design and Implementation:**
   o **Tool Selection:** Choose either Power BI or RapidMiner for the ETL implementation. Justify your choice.
   o **Data Sources:** Describe the simulated or actual data sources you will use for your implementation. If using simulated data, explain how you generated it. If using real data, ensure you have the necessary permissions and anonymize any sensitive information.
   o **ETL Steps:** Detail the steps involved in your ETL process:
      ▪ **Extraction:** Describe how you will extract data from the various sources.
      ▪ **Transformation:** Explain and implement the transformations you will perform:

- Data Cleaning: Handle missing values (using appropriate techniques like mean/median imputation, mode imputation for categorical data), remove duplicates, correct inconsistencies. Justify your chosen methods.
- Data Type Conversion: Ensure proper data types.
- Calculations: Derive new measures (e.g., total sales, average interaction time).
- Data Integration: Join data from different sources.
- Surrogate Key Generation: Create surrogate keys for dimension tables.
- Date Dimension Creation: Create a date dimension table.
- Data Transformation: Include log transformation for skewed numerical data (if applicable) and one-hot encoding for categorical variables.
  - **Loading:** Describe how you will load the transformed data into the data warehouse (Power BI data model or RapidMiner repository).
4. **Data Warehouse Implementation (Power BI or RapidMiner):**
   o Implement the designed schemas in your chosen tool.
   o Establish relationships between fact tables and dimension tables.
5. **Analysis and Reporting (Power BI or RapidMiner):**
   o Develop at least three reports or visualizations that demonstrate the analytical capabilities of your data warehouse. These should showcase insights related to the chosen business processes and, ideally, cross-process analysis enabled by the fact constellation schema.

**Power BI Tool:**

Power BI is a Data Visualization and Business Intelligence tool that converts data from different data sources to interactive dashboards and BI reports. In Power BI we can load the data from sources like csv, excel files and also can transform that data. In here we can easily design different styles of data mart schemas which are most widely used to develop data warehouses.

**Dataset Used:**

For this assignment we have used one dataset. It is the **Supermarket sales dataset**. We have downloaded these datasets from the Kaggle.

**Loading Supermarket Sales Dataset In Power**

After performing preprocessing on the dataset we moved towards the next step which is the "Loading" step. We have opened the Power BI tool and there is the Get Data option present in the Home tab. Here we have given the Supermarket sales dataset file that we have to load and after giving the file path it shows a preview of the dataset in the window and then click on the Load button to load the dataset into Power BI. We successfully loaded the Dataset in Power BI Tool.

**Star Schema:**

The Star schema is the simplest style of data fomat schema which is most widely used to develop data warehouses. It consists of a fact table in the middle of the schema connecting to a set of dimension tables. The fact table and dimension table hold the columns that store the data for the model.

Fact table mainly consists of the business facts and foreign keys that refer to the primary keys in the dimension table. The Dimension table consists mainly of descriptive attributes.

**Creating Star Schema:**

When we have loaded the data then there exists only single table names as Supermarket sales tables which act as our FactTable. We have renamed the table as Fact Supermarket Sales.

So firstly create a dimension table named **Dim Customer** for the Customers where we have attributes Customer_id, Gender and Customer_Type . After doing that we performed a merge operation on column Customer_id from Supermarket Sales fact table and newly created Customer dimension table and this is how we added a foreign key Dim Customer.Customer_id to Supermarket fact table and deleted the rows Gender and Customer_Type from this table.

Similarly we created a dimension table for date named **Dim Date-time** where we added the columns Date and Invoice id, Time to new table and removed the duplicates and given ids to it and merged it with our main table and then deleted the columns from Supermarket Sales fact table



Similarly we did for Branch and City column and created a new dimension table named **Dim Location** with Branch, City and Location_id.

Similarly for payment attributes we have created a new dimension table named **Dim Payment** method with payment and Payment_method_ID as primary key.



We created **Dim Product**, with Product line, Product_id and Unit price as attributes.

We created **Dim Profits**, with gross income and Invoice id as attributes.



That completes our Star Schema.

**Snowflake Schema:**

Snowflake Schema is basically a variation of Star Schema which is represented by centralized Fact Tables which are connected to multiple Dimension Tables. But the difference between Star schema and Snowflake schema is that in Snowflake schema large dimension tables are normalized into multiple tables.

**Creating Snowflake Schema:**

For creating Snowflake schema we have copied Invoice Id attributes from the Dim Profits table into another table named as Dim Date-Time table. After this we have removed duplicates from attributes and created a new column Invoice ID that acts as the primary key into the table and linked the table with the Dim Date-Time table. That completes our Snowflake Schema.
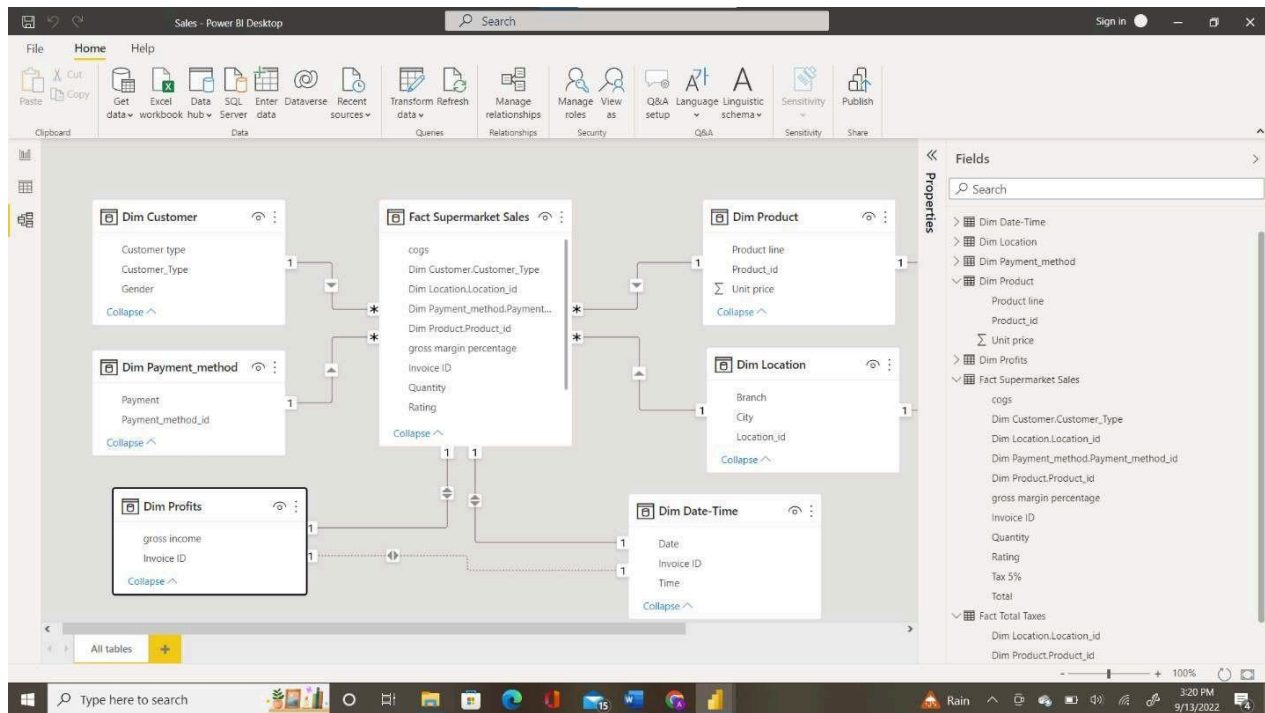
**Fact Constellation Schema:**

Fact Constellation is a schema for representing a multidimensional model. It is a collection of multiple fact tables having some common dimension tables. It can be viewed as a collection of several star schemas and hence, also known as **Galaxy schema.**
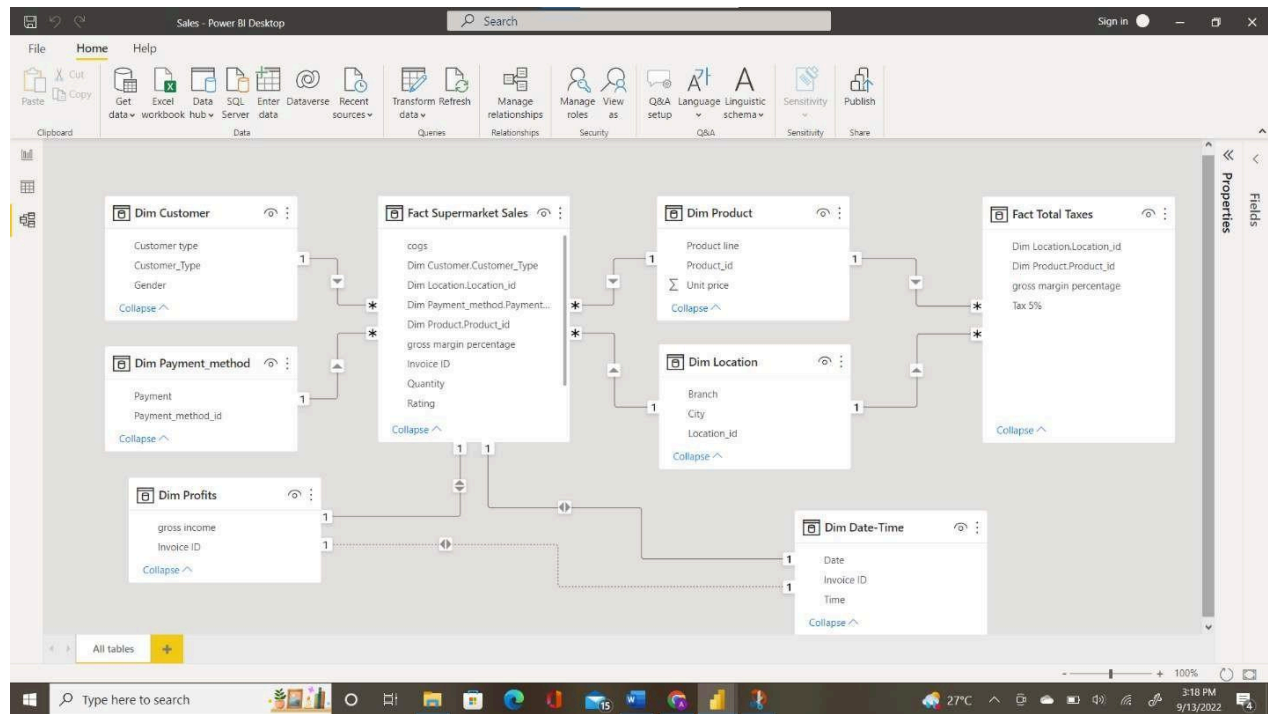
**Creating Fact Constellation Schema:**

For creating Fact Constellation Schema we have created a new Fact Total taxes, which acts as another Fact Table for our schema. In the Dim Product dimension table which is linked to Fact Supermarket Sales Table there is one Product id attribute which is also present in Fact Total Taxes Table, hence we have linked the Dim Product table to the Fact Total Taxes Table also. That means the Dim Product table is now linked to both Fact Tables.

Similarly In the Dim Location dimension table which is linked to Fact Supermarket Sales Table

there is one Location id attribute which is also present in Fact Total Taxes Table, hence we have linked the Dim Location table to the Fact Total Taxes Table also. That means the Dim Location table is now linked to both Fact Tables.

**Conclusion:**

Thus, we have successfully preprocessed and loaded the dataset into the Power BI tool and designed Star, Snowflake and Fact Constellation Schemas.

**RESULTS**

### 1. Business Process Chosen:

**Health and Lifestyle Analysis in the Food Industry** – focusing on analyzing food preferences, eating habits, lifestyle choices, and their impact on student well-being.

### 2. Schema Design:

Designed a **Star Schema** with the central **Fact Table** capturing student food behavior metrics (e.g., calories, meals out, nutritional checks) and **Dimension Tables** including student demographics, eating habits, cuisine preferences, family background, and health indicators.

### 3. ETL Process:

I. **Extracted** data from a publicly available dataset using Power BI's **Power Query Editor.**

II. **Transformed** data by handling nulls, correcting data types (e.g., GPA), removing unnecessary columns, and filtering irrelevant rows (e.g., from row 125 onwards).

III. **Loaded** the clean data into Power BI's model and established relationships across the tables to create an optimized data warehouse structure.

### 4. Tools Used:

I. Power BI (Data loading, transformation, schema design, and visualization)

II. Power Query Editor (for ETL operations)

**Results Attached:**

## I.  Data Loading in Power BI



**Table : Student_Info**



**Table : Eating_Habits**

## II.  Star Schema Diagram (Model View)



**Star Schema**

## III.  Creating new measure using Power Query



```
1  AvgCaloriesPerDay = AVERAGE('Food_and_Calories'[calories_day])
```



```
1  EmployedByIncome =
2  CALCULATE(
3      COUNTROWS(Student),
4      Student[employment] <> BLANK()
5  )
6
```