Sri Lanka Institute of Information Technology

# Assignment I

Data Warehouse & Business Intelligence

2022

Submitted by:

Mayurrsh.T

Y3S1.04 (DS)

# Contents

# 1.Data set selection

Provided by: kaggle.com
Data Set Name: Melbourne Housing Snapshot
Data Set:
https://www.kaggle.com/datasets/dansbecker/melbourne-housing-snapshot

About Dataset

The chosen data source is a Kaggle based collection of transactional data. Which represents Melbourne's house-sale information. It is made up of a single CSV file containing enough data in 21 columns. The original huge CSV file has been divided into smaller sub-CSV files. like Seller Details and Property, New Identifiers are contained in the sub-CSV files. In addition, I manually changed some data records to meet the requirements.
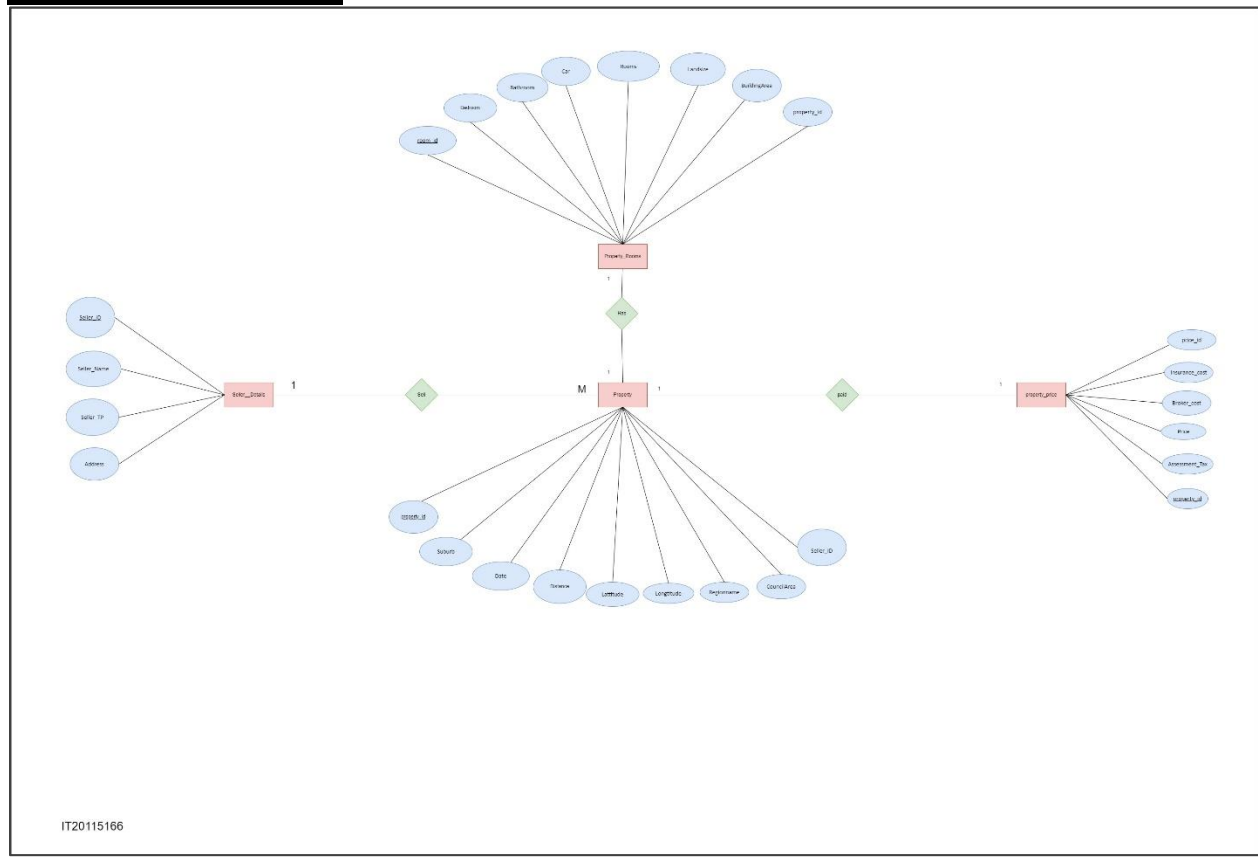
## 2.PREPARATION OF DATA SOURCE

All the data sources are provided in csv format by the web site. In preparation of data sources, some changes have done for the source format. Some of the given files were converted into text files and Property Details csv files into a source database, while others were removed and added to another file.

Final State of Preparation of the source data formats before Transforming data:
- ➢ Property_Details.csv
- ➢ Property_Price.csv
- ➢ Property_Rooms.csv
- ➢ Seller.txt

This data collection contains information about house sales in Melbourne. It comprises information on the residence as well as information about the sellers who sold such properties. One csv file was split into two halves, as shown below.

# ER-Diagram



IT20115166

- ➢ The above diagram shows the connection between the entities in the data set
- ➢ Assumptions:
  - One Seller have many Properties.
  - There can be many campaign data sets in a single summary report
  - Many client transactions are summarized in a single summary report.

# Description of the Data Set

| Source Type | Table Name | Include |
|---|---|---|
| Seller.txt | Seller | <table><thead><tr><th>Column</th><th>Data type</th><th>Description</th></tr></thead><tbody><tr><td>Seller_ID</td><td>nvarchar(255)</td><td>Unique id of Seller (PK)</td></tr><tr><td>Seller_Name</td><td>nvarchar(255)</td><td>Name of Seller</td></tr><tr><td>Seller_TP</td><td>nvarchar(255)</td><td>Phone number of Seller</td></tr><tr><td>Address</td><td>nvarchar(255)</td><td>Address of Seller</td></tr></tbody></table> |
| Melbourne_Housing_SnapshotDB | Property_Details | <table><thead><tr><th>Column</th><th>Data type</th><th>Description</th></tr></thead><tbody><tr><td>property_id</td><td>int</td><td>Unique id of Property (PK)</td></tr><tr><td>Suburb</td><td>nvarchar(255)</td><td>Name of Residential name</td></tr><tr><td>Address</td><td>nvarchar(255)</td><td>Adress of Property</td></tr><tr><td>Method</td><td>nvarchar(255)</td><td>Sold Method</td></tr><tr><td>Distance</td><td>float</td><td>Distance From Capital</td></tr><tr><td>Seller_ID</td><td>nvarchar(255)</td><td>Unique id of Seller (FK)</td></tr><tr><td>Lattitude</td><td>float</td><td>Lattitude</td></tr><tr><td>Longtitude</td><td>float</td><td>Longitude</td></tr><tr><td>Regionname</td><td>nvarchar(255</td><td>Name of Regional</td></tr><tr><td>CouncilArea</td><td>nvarchar(255</td><td>Governing Council for the Area</td></tr></tbody></table> |
| | Property_Price | <table><thead><tr><th>Column</th><th>Data type</th><th>Description</th></tr></thead><tbody><tr><td>price_id</td><td>int</td><td>Unique id of Price (PK)</td></tr><tr><td>Assessment_Tax</td><td>float</td><td>Price of Aessment Tax</td></tr><tr><td>Broker_cost</td><td>float</td><td>Price of Broker Cost</td></tr><tr><td>Insurance_cost</td><td>float</td><td>Price of Insurance Cost</td></tr><tr><td>Price</td><td>float</td><td>Price of that Property</td></tr><tr><td>property_id</td><td>int</td><td>Unique id of Property (FK)</td></tr></tbody></table> |
| | Room_Count | <table><thead><tr><th>Column</th><th>Data type</th><th>Description</th></tr></thead><tbody><tr><td>room_id</td><td>int</td><td>Unique id of Room (PK)</td></tr><tr><td>Bedroom</td><td>int</td><td>Number of bedrooms in this Property</td></tr><tr><td>Bathroom</td><td>int</td><td>Number of bathrooms in this Property</td></tr><tr><td>Car</td><td>int</td><td>Number of Cars spots in this Property</td></tr><tr><td>Rooms</td><td>int</td><td>Number of rooms in this Property</td></tr><tr><td>property_id</td><td>int</td><td>Unique id of Property (FK)</td></tr><tr><td>Landsize</td><td>int</td><td>Size of the Land</td></tr><tr><td>BuildingArea</td><td>int</td><td>Area of the Building</td></tr></tbody></table> |

# Design of Data_Source

- **Property_Details**



- **Property_Price**

- **Property_Rooms**

# 3.SOLUTION ARCHITECTURE



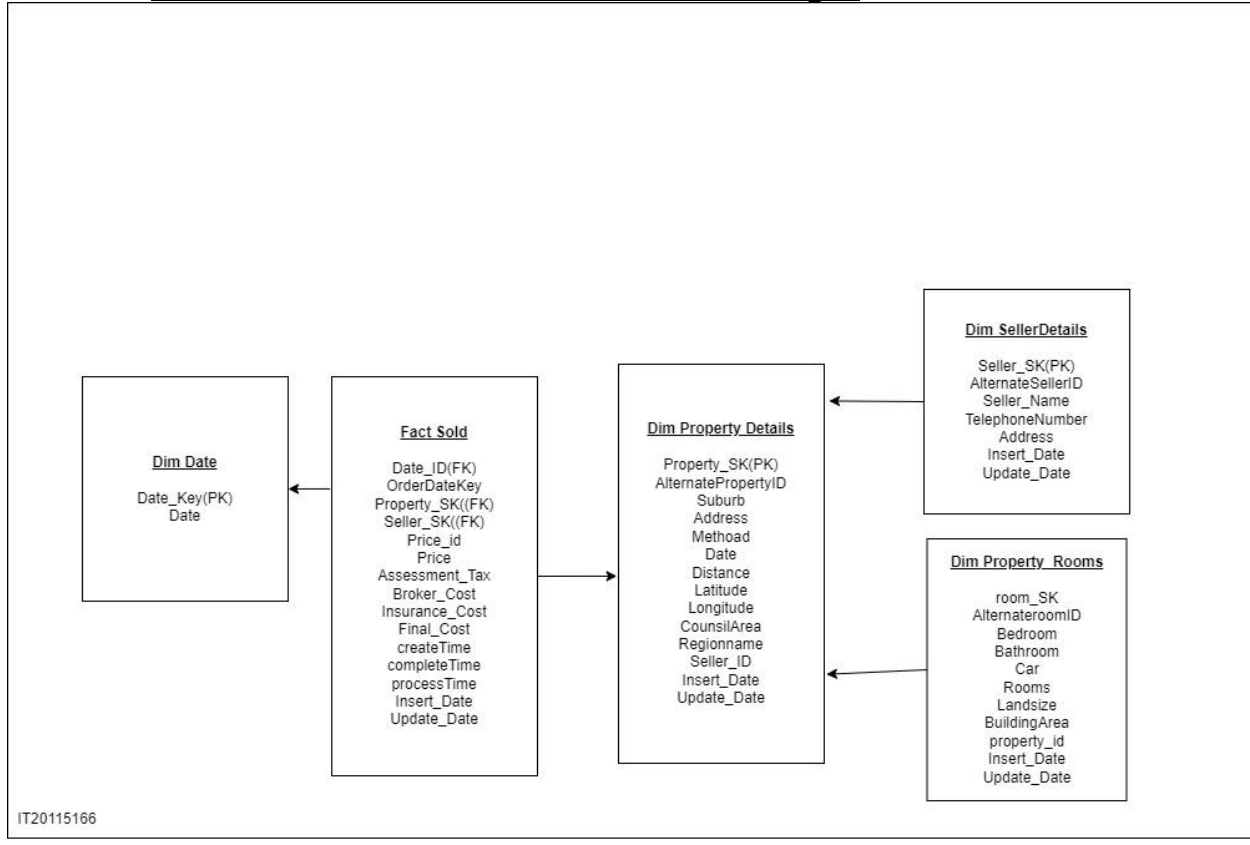(As the figure 2 shows for the ETL processing, initially)
- **Seller**: Text File
- **Melbourne _Housing _SnapshotDB**: Source Database,

We can handle data from various sources and transform it to business insights to make decisions, analyze data, and produce reports using diverse procedures, structures, and technologies. This will also give the data a new dimension

# 4.Data warehouse Design & Development

## I.    Data warehouse Table design



| Dimention Name | Dimention Attributes | Data Type | Key column | Derived Logic |
|---|---|---|---|---|
| Dim_Property_ Details | Property_SK | int | Primary key | Auto increment |
| | AlternateProper tyID | int | | |
| | Suburb | varchar (50) | | |
| | Address | varchar (50) | | |
| | Method | varchar (50) | | |
| | Date | varchar (50) | | |
| | Diatance | varchar (50) | | |

| | | | | |
|---|---|---|---|---|
| | Lattitude | varchar (50) | | |
| | Longitude | varchar (50) | | |
| | Seller_ID | int | | |
| | Regionname | varchar (50) | | |
| | Insert_Date | datetime | | System Datetime |
| | ModifiedDatedate | datetime | | System Datetime |
| | | | | |
| | | | | |
| Dim_SellerDetails | Seller_SK | int | Primary key | |
| | AlternateSellerID | int | | |
| | Seller_Name | nvarchar(50) | | |
| | Telephone_Number | nvarchar(50) | | |
| | Address | nvarchar(50) | | |
| | Insert_Date | datetime | | System Datetime |
| | Update_Date | datetime | | System Datetime |
| | | | | |
| DimDate | DateKey | int | Primary key | |
| | Date | datetime | | |
| | FullDateUK | char(10) | | |
| | FullDateUSA | char(10) | | |
| | DayOfMonth | varchar (4) | | |
| | DaySuffix | varchar (9) | | |
| | DayName | varchar (9) | | |
| | More…. | | | |
| | | | | |
| Fact_Sold | Seller_SK | int | foreign key | |
| | Property_SK | int | foreign | |

| | | | key | |
|---|---|---|---|---|
| | Date_ID | int | foreign key | |
| | Price | float | | |
| | price_id | int | | |
| | Assessment_Tax | varchar (50) | | |
| | Broker_Cost | varchar (50) | | |
| | createTime | datetime | | |
| | completeTime | datetime | | |
| | processTime | datetime | | |
| | Insurance_Cost | float | | |
| | Final_Cost | float | | Price+Assessment_Tax+Broker_Cost+Insuarence_Cost |
| | Insert_Date | datetime | | System Datetime |
| | Update_Date | datetime | | System Datetime |
| | | | | |
| DimProperty_Rooms | room_SK | int | | |
| | AlternateroomID | int | | |
| | Bedroom | varchar (50) | | |
| | Bathroom | varchar (50 | | |
| | Car | varchar (50 | | |
| | Rooms | varchar (50 | | |
| | landsize | varchar (50 | | |
| | BuildingArea | varchar (50 | | |
| | property_id | int | | |
| | Insert_Date | datetime | | System Datetime |
| | Insert_Date | datetime | | System Datetime |

# Calculation

: Final Cost =( Price+ Assessment_Tax+ Broker_Cost+ Insurance_Cost)

# I. Assumptions

- dbo.DimDate is added to the Data Warehouse for better performance.
- dbo. Property_Price is used in creating the fact table

# II. Slowly changing dimensions

- Customer Details were considered as a slowly changing dimension

| • Dimension table | Attributes |
|---|---|
| Dim_SellerDetails | Telephone_Number (changing attribute) Address (Hostorical) |

# 5.ETL Development

## I. Data Extraction & Load into Staging Tables
## Property_Details



(Property Details is extracted from Property Details the table in the source database and inserted to the Property Details Staging table)

## Property_Rooms



(Place Details is extracted from Property Rooms  the table in the source database and inserted to the  Property Place Staging table)

# Property_Price



(Price Details is extracted from Property_Price the table in the source database and inserted to the Property Price Staging table)

# Property Rooms



(Room Details is extracted from room_count the table in the source database and inserted to the Room Count Staging table)

# Seller_Details



(Seller Details is extracted from Seller.txt the table in the source database and inserted to the Seller Staging table)



(The Control Flow of 'Extract Data and Load into Staging' Step can illustrate as the give figure)

(Staging Tables created and values inserted)



(Available are In Staging

## II.    Data Profiling

Data Profiling provides the means of analyzing large amount of data using different kind of processes. In this step, null values, repeated values and quality of the data is checked.



- ❖ Each staging table is profiled and saved in a specific folder.
- ❖ As the figure shows, after the Staging step doing this task shows the things what the developer has to consider about the data which are stored in staging table and the developer is able to identify the issues with staging data by data profiling (such as null values).
- ❖ The diagram depicted the entirety of Data Profiling as it relates to Staging.

## III.   <u>Data Transformation and Loading</u>

- Data Transformation is developed according to the dimensional modeling designed above.



- In this step, the Dimension Tables created in Melbourne Housing Snapshot_DW are loaded with the data of relevant staging tables

(DimProperty_Rooms)



❖ Property Room data is loaded to the DimProperty_Room.

❖ Sort and merge transformation tasks are used.



❖ Update DimProperty_Room procedure is used to check whether the data inserted or not.

# (DimProperty_Details)



❖ Property Details data is loaded to the DimProperty_Details.

❖ Sort and merge transformation tasks are used.



❖ Update DimProperty_Details procedure is used to check whether the data inserted or not
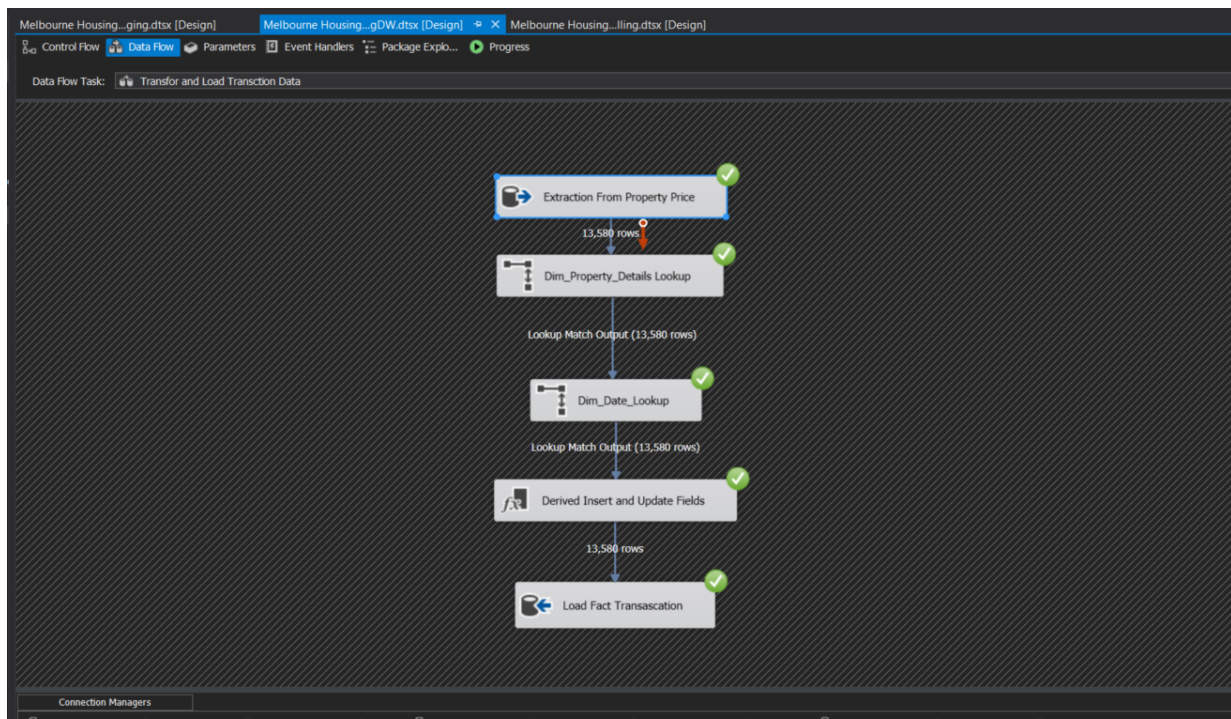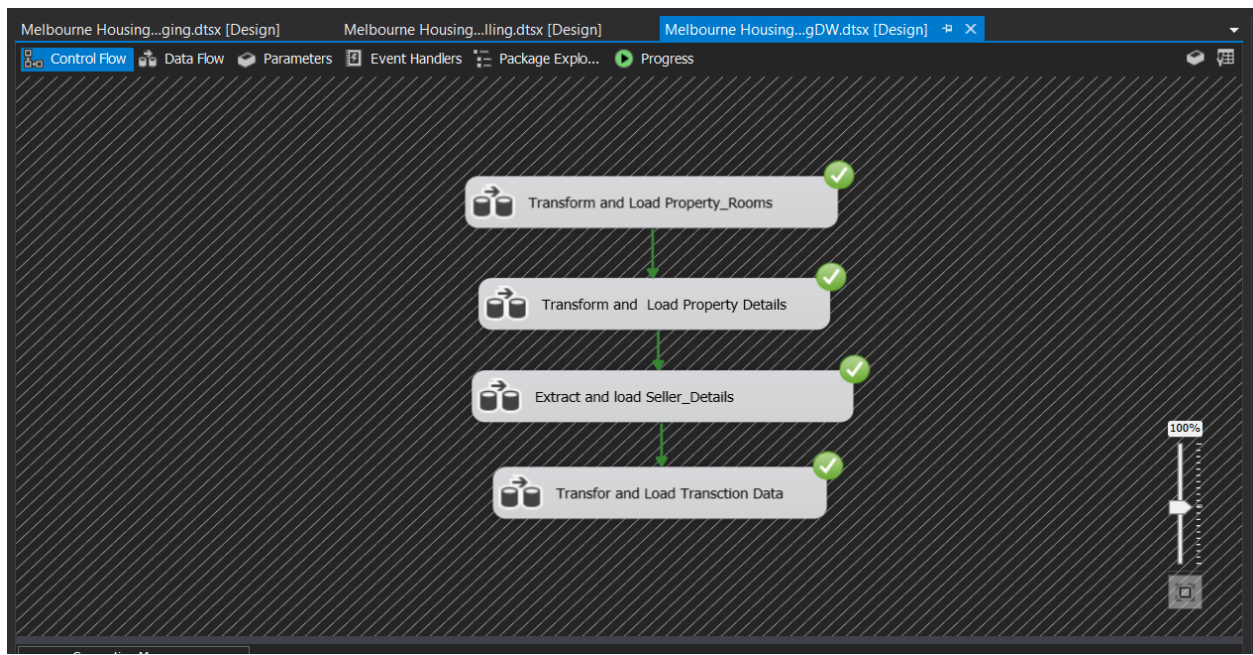
# Loading Slowly Changing Dimension

- ✛ DimSellerDetails is the slowly changing dimension in this dimensional modeling.
- ✛ In order to load data to Dimension table, the slowly changing dimensions (historical) have two specific columns as StartDate & End Date to ensure that the data is valid at the moment.
- ✛ slowly changing dimension wizard let the developer to select the Dimension table, Business keys of the dimension and what would be the slowly changing attributes.
- ✛ The below mentioned columns were set as changing attributes:
  - o Telephone_Number : Changing Attribute
  - o Adresss :Histrorical Attribute

# Load data to Fact table

- The final step of Transformation & Loading is load data to fact table. According to the dimensional model, TransactionStaging table is used to insert values into DimTransaction table.

- After loading to all the dimensions, lastly data was loaded to the fact table. The below steps were followed:

    - ❖ Data extracted from the StgProperty_Price staging.

    - ❖ Join operation is done for the Dim_Property_Details Lookup.

    - ❖ Join operation is done for the Dim_Date_Lookup.

    - ❖ insert and modified date were derived.

    - ❖ Fact details loaded to the Fact_Load table.

All Details are loaded