



# Real Estate Sales

## 2001–2022

"A Visual Exploration of Property Sales  
Trends, Locations, and Growth"

This project explores over two decades of real estate data to uncover patterns, trends, and insights across various towns and property types. Using interactive dashboards, the analysis helps identify high-performing areas, pricing shifts, and key market indicators. The aim is to provide a clear and visual understanding of how the property market has evolved between 2001 and 2022, making it useful for both data enthusiasts and real estate professionals.

Prepared by: Anurag Yadav

Dated: 25-07-2025

## 2. Abstract

This project presents a comprehensive and interactive real estate sales dashboard built using Power BI, enhanced with a dynamic flow animation created in Python. The primary objective is to analyze and visualize real estate data effectively while ensuring data quality and user-friendly navigation.

The dashboard is divided into logical sections—**Overview**, **Sales Analysis**, **Property Details**, **Remarks & Insights**, and **Data Quality**—to help stakeholders quickly grasp key trends, identify outliers, and assess the integrity of the data. Python was used extensively in the data preparation phase to clean, transform, and identify missing or anomalous values. A unique feature of this project is the animated flow diagram developed using matplotlib, which illustrates the logical transition between dashboard pages. This visual serves as both an intuitive guide and a project explainer.

By combining Python's flexibility in data handling and animation with Power BI's powerful visualization capabilities, the project highlights a practical, visually engaging approach to data storytelling in real estate analytics.



## Table of Contents

<b>1. Cover Page .....</b>	<b>1</b>
<b>2. Abstract .....</b>	<b>1</b>
<b>3. Table of Contents .....</b>	<b>2</b>
<b>4. Introduction.....</b>	<b>4</b>
4.1 Background .....	4
4.2 Project Overview.....	4
4.3 Project Objectives .....	5
<b>5. Data Description .....</b>	<b>5</b>
5.1 Data Source & Acquisition .....	5
5.2 Data Characteristics .....	6
<b>6. Methodology .....</b>	<b>7</b>
6.1 Data Cleaning & Preprocessing (Python) .....	7
6.2 Database Design & Querying (MySQL/SQL).....	8
6.3 Data Visualization & Dashboarding (PowerBI).....	8
6.4 Outcome .....	9
<b>7. Analyzing Real Estate Market Trends for Better Understanding .....</b>	<b>9</b>
7.1 Datasets explored in this project .....	9
7.2 Business Impact & Insight Potential .....	10
<b>8. Python Programming for Data Cleaning &amp; ETL .....</b>	<b>12</b>
8.1 ETL Process Flow Diagram.....	12
8.2 Directory Hierarchy.....	13
8.3 Module Description .....	13
8.3.1 main.py .....	13
8.3.2 splitting.py .....	14
8.3.3 cleaning.py .....	15
8.3.4 datatype_utils.py .....	16
8.3.5 dateconversion.py .....	17
8.3.6 read_dfs.py .....	18
8.3.7 database_connection.py.....	19
8.3.8 sql_generator.py .....	20

<b>9. Conclusion &amp; Future Scope .....</b>	<b>21</b>
9.1 Achievements.....	21
9.2 Scope for future improvement .....	21
9.2.1 Error logging & monitoring.....	21
9.2.2 Unit Testing.....	21
9.2.3 Configure file Integration.....	21
9.2.4 Support for other Data source.....	21
9.2.5 Parallel Processing .....	21
9.2.6 Integration with orchestration tools.....	22
9.2.7 Data Validation Layer.....	22
<b>10. Uncovering Trends in the Real Estate Market Using PowerBI .....</b>	<b>23</b>
10.1 Introduction .....	23
10.2 Dashboard Structure.....	23
10.3 Main Dashboard Pages .....	24
10.3.1 Overview .....	24
10.3.2 Sales Analysis .....	25
10.3.3 Property Details .....	26
10.3.4 Remarks & Insights .....	27
10.3.5 Data Quality .....	28
10.4 Drill through Pages.....	29
10.4.1 Property Town stats.....	29
10.4.2 Residential Property stats.....	30
10.5 Interactive Features .....	31
10.6 Value Proposition.....	31
<b>11. Discussion .....</b>	<b>32</b>
<b>12. Limitation .....</b>	<b>33</b>
<b>13. Future Scope .....</b>	<b>33</b>
<b>14. Conclusion .....</b>	<b>34</b>
<b>15. Analyst Profile .....</b>	<b>34</b>

## 4. Introduction

### 4.1 Background

The real estate industry generates vast amounts of data ranging from property transactions and sales trends to location-specific attributes and customer behavior. Efficiently analyzing and interpreting this data is crucial for informed decision-making by realtors, investors, buyers, and analysts. However, the challenge often lies in organizing this information in a meaningful way, ensuring its quality, and presenting it in an accessible format for various stakeholders.

### 4.2 Project Overview

This project focuses on developing a real estate sales dashboard using Power BI, integrated with custom Python-based animations to guide the user journey. The dashboard is structured into five major sections:

- **Overview**
- **Sales Analysis**
- **Property Details**
- **Remarks & Insights**
- **Data Quality**

The data used in the dashboard undergoes cleaning and transformation in Python before being imported into Power BI. Additionally, a unique flow diagram was created using matplotlib to visualize the navigation and logical flow between dashboard sections. This animated diagram acts as a landing page element, offering a high-level summary of the dashboard's structure.

## 4.3 Project Objectives

The main objectives of this project are:

- To analyze real estate data and uncover actionable insights related to property sales and attributes.
- To assess data quality through null value detection, outlier identification, and column-level integrity checks.
- To create an intuitive and visually compelling dashboard using Power BI.
- To enhance user navigation by integrating a custom-built Python flow animation diagram on the landing page.

By bridging Python's data processing and visualization capabilities with Power BI's interactive dashboards, this project aims to provide a professional-grade analytical tool for real estate performance monitoring and decision support.

## 5. Data Description

### 5.1 Data Source & Acquisition

The dataset used in this project was sourced from [Kaggle.com](#), a well-known platform for data science and machine learning datasets. The dataset, titled "**Real Estate Sales 2001–2022**", was uploaded by **Omnia Mahmoud Saeed**. It provides detailed historical data on property sales recorded between the years 2001 and 2022.

The dataset was downloaded in CSV format and initially processed using Python (pandas and matplotlib) for cleaning, transformation, and exploratory data analysis. The cleaned data was then stored in a relational database (MySQL) and connected to **Power BI** for building interactive visual dashboards.

## 5.2 Dataset Characteristics

This dataset offers a comprehensive view of real estate transactions with a mix of numerical, categorical, and geographic features. It includes information about property characteristics, financial assessments, sale outcomes, and temporal details. Below is an overview of the key columns:

Column Name	Description
Serial Number	A unique identifier for each property record.
List Year	The year when the property was listed for sale.
Date Recorded	The exact date the property sale was recorded.
Town	The town or city where the property is located.
Address	Street-level address of the property.
Assessed Value	Value assigned for property tax purposes.
Sale Amount	Final transaction amount at which the property was sold.
Sales Ratio	Ratio of the sale amount to the assessed value—indicating valuation accuracy.
Property Type	General type of property (e.g., Residential, Commercial).
Residential Type	More specific residential classification (e.g., Single Family, Condo).
Non Use Code	Flags special-use or vacant properties.
Assessor Remarks	Notes or observations made by property assessors.
OPM Remarks	Comments from the Office of Property Management.
Location	Geographic coordinates (latitude and longitude) of the property.

In addition to the columns listed above, the dataset also includes:

- Price Ranges:** Binned price categories for better segmentation and comparative analysis.
- Time Period Buckets:** Grouped sales by year or month for trend analysis.

- **Geospatial Tags:** Latitude and longitude values for spatial visualization and mapping.

This diverse and rich dataset enables comprehensive real estate market analysis by providing both structured numeric fields and contextual descriptive fields across different regions and time periods.

## 6. Methodology

This project followed a structured ETL (Extract, Transform, Load) pipeline to ensure that raw real estate sales data was cleaned, stored, analyzed, and visualized effectively. The methodology is broken into four key phases:

### 6.1 Data Cleaning & Pre-Processing (Python)

Data cleaning and preprocessing were performed using Python libraries such as *pandas*, *numpy*, and *matplotlib*. The main tasks included:

- **Handling Missing Values:** Columns with null or empty entries were identified and treated using appropriate strategies like deletion, imputation, or tagging.
- **Standardizing Formats:** Date fields were converted into datetime format, and text fields were stripped of unnecessary whitespace and special characters.
- **Feature Engineering:** New columns were derived such as price bins, year buckets, and categorized property types to aid in grouping and filtering.
- **Latitude/Longitude Separation:** The Location column was split into separate Latitude and Longitude columns for mapping.
- **Outlier Detection:** Box plots and scatter plots were used to identify and optionally remove extreme outliers in Sale Amount and Sales Ratio.
- **Export to Database:** The cleaned data was saved into structured CSVs and then loaded into a MySQL database for further use.

## 6.2 Database Design & Querying (MySQL/SQL)

The cleaned dataset was imported into a MySQL database, where data modeling and querying took place:

- **Table Structure:** Separate tables were created for Property\_Sales, Time\_Reference, and Geographic\_Data to normalize the dataset.
- **Primary/Foreign Keys:** Unique identifiers such as Serial Number were used to relate tables where appropriate.
- **SQL Queries:** Multiple queries were written to:
  - ✓ Aggregate total and average sale amounts by town or year
  - ✓ Count missing values per column
  - ✓ Calculate overall sales ratio and price ranges
  - ✓ Filter property types dynamically based on sales trends
- **Views:** Pre-aggregated views were created for optimized Power BI consumption.

## 6.3 Data Visualization & Dashboarding (Power BI)

Power BI was used to transform data insights into interactive visual dashboards:

- **Multi-Page Dashboard:** Organized into sections like Overview, Trends, Property Types, and Outliers.
- **Key Visuals:**
  - ✓ Column & Bar Charts for sales comparisons
  - ✓ Scatter Plots for outlier detection
  - ✓ Word Clouds for textual insights from remarks
  - ✓ Map Visuals using latitude/longitude
  - ✓ Python Visuals to show missing values
- **Slicers & Filters:** Dynamic controls were added for towns, years, property types, and price ranges.
- **Navigation Buttons:** Custom buttons allow seamless navigation between pages using bookmarks and drillthrough functionality.

- **Design Theme:** A consistent and visually appealing theme was applied using contrasting colors and minimalist design.

## 6.4 Outcome

The final outcome of this project is an ***insight-rich Power BI dashboard*** that enables users to:

- Identify sales trends over time and across towns
- Compare sale amounts against assessed values
- Detect outliers and pricing anomalies
- Filter and drill through different property types and years
- View geospatial distributions of properties
- Assess data quality by visualizing missing values

The project offers a complete workflow from raw data to business-ready insights, which can assist stakeholders like property analysts, investors, or urban planners in making data-driven decisions.

## 7. Analyzing Real Estate Market Trends for Better Understanding

This section focuses on how the data was explored and analyzed to uncover key market dynamics, patterns, and opportunities within the real estate sector. By leveraging advanced data exploration techniques, the project uncovers critical insights that can benefit both business users and analysts.

### 7.1 Datasets Explored in This Project

In this project, the "**Real Estate Sales 2001–2022**" dataset from **Kaggle** (uploaded by *Omnia Mahmoud Saeed*) was the primary data source. The dataset provides comprehensive transactional and assessment information across multiple years and towns. The following key datasets (or structured tables after cleaning) were explored:

#### ✓ **Property Sales Table**

Includes columns like Serial Number, Date Recorded, Town, Address, Sale Amount, Assessed Value, Sales Ratio, and Property Type.

#### ✓ **Categorical Groupings**

Price ranges, time periods, and grouped property types were engineered to aid in comparative analysis and segmentation.

#### ✓ **Geolocation Data**

Latitude and longitude from the Location column were parsed to support map-based visualizations.

#### ✓ **Remarks & Descriptions**

Text columns such as Assessor Remarks and OPM Remarks were used to generate insights via word clouds.

By thoroughly exploring these datasets through filtering, aggregation, and visualization, a robust understanding of the real estate landscape was formed.

## 7.2 Business Impact & Insight Potential

The insights derived from this project have strong potential to support a range of real estate and urban planning decisions. Some of the major business impacts include:

- **Pricing Strategy Optimization**

- ✓ By analyzing the Sales Ratio, users can identify areas where properties are undervalued or overpriced.
- ✓ Trends in assessed vs. sale prices can guide appraisal strategies and fair value assessments.

- **Market Timing & Investment Decisions**

- ✓ Time series analysis reveals high and low seasons for property sales.
- ✓ Helps real estate investors decide *when* and *where* to buy or sell.

- **Urban Development & Policy Making**

- ✓ Aggregated town-wise data enables urban planners to identify regions of rapid development or stagnation.
- ✓ Supports zoning decisions based on property type distribution.

- **Data Quality Audits**
  - ✓ Visualizing missing values helps data stewards or analysts focus on areas where improvements in record-keeping or system integration are needed.
- **Outlier Detection for Fraud Prevention**
  - ✓ Sudden spikes or anomalies in sale prices are easily spotted through scatter plots, alerting to potential fraud or misreporting.

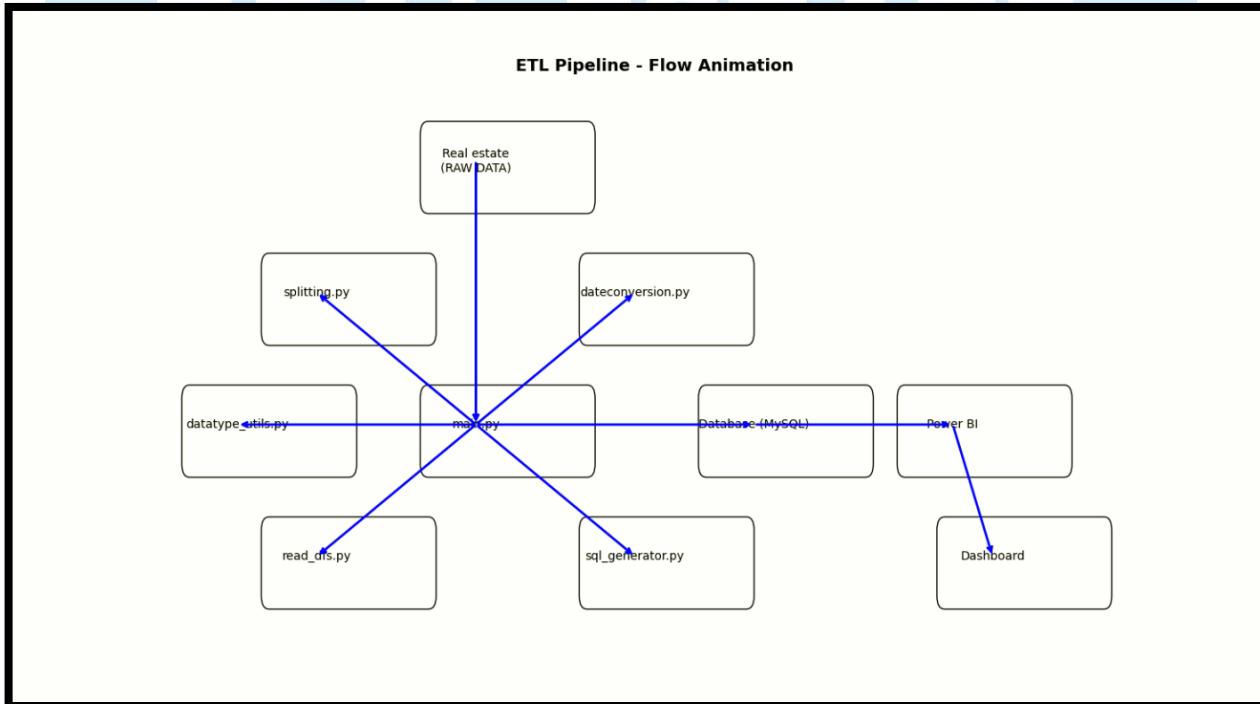
Overall, this project bridges data engineering, analytics, and visualization to support actionable business decisions within the real estate domain.



## 8. Python Programming for Data Cleaning & ETL

In this project, Python is used to automate the entire data engineering workflow, starting from *splitting the raw dataset*, followed by *cleaning, transforming*, and *loading the data into a MySQL database*. The structure of the program is modular to enhance reusability and scalability.

### 8.1 ETL Process Flow Diagram



**Figure 1:** Flow diagram showing how `main.py` orchestrates the full ETL pipeline by invoking modules such as `splitting.py`, `cleaning.py`, `datatype_utils.py`, and ultimately `sql_generator.py` to create and populate database tables.

## 8.2 Directory Hierarchy

```
Python_Programme/
|
└── main.py                                # Main execution script
|
└── etl_utils/
    ├── splitting.py                         # Sub-directory containing utility scripts
    ├── cleaning.py                          # Splits raw files into smaller CSVs
    ├── datatype_utils.py                   # Cleans individual dataframes
    ├── dateconversion.py                 # Converts datatypes for DB compatibility
    ├── read_dfs.py                         # Standardizes date columns
    ├── database_connection.py            # Reads splitted CSVs into pandas DataFrames
    └── sql_generator.py                  # Establishes MySQL DB connection
                                         # Generates SQL tables & inserts data
```

## 8.3 Module Description

### 8.3.1 main.py

- Purpose: This is the central execution script which coordinates the entire ETL process.
- Description:
  - ✓ Imports utility functions
  - ✓ Initiates splitting, cleaning, and DB insertion
  - ✓ Manages process flow in a structured manner
- Code Snippet:

```
from etl_utils import splitting, cleaning, read_dfs, sql_generator

def main():
    splitting.split_file("raw_dataset.csv")
    dfs = read_dfs.load_all_csvs()
    cleaned_dfs = [cleaning.clean(df) for df in dfs]
    sql_generator.insert_into_db(cleaned_dfs)

if __name__ == "__main__":
    main()
```

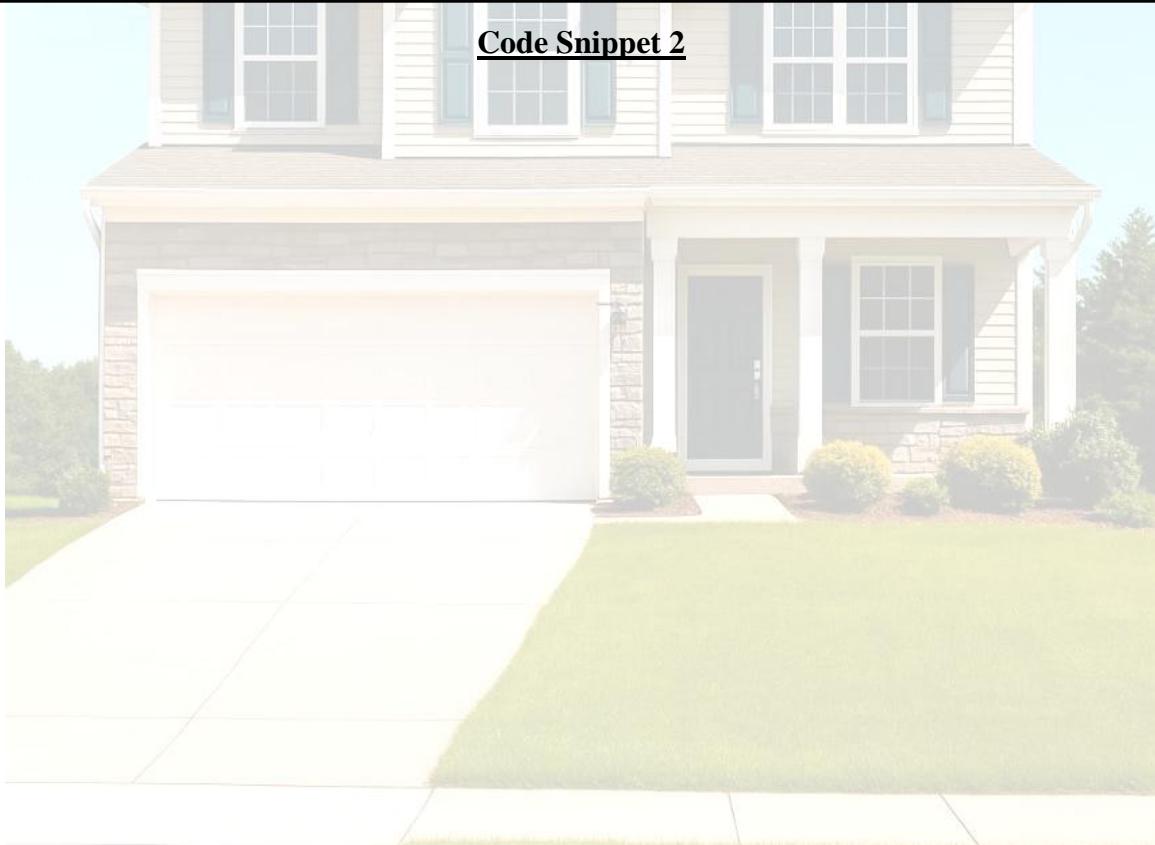
**Code Snippet 1**

### 8.3.2 splitting.py

- Purpose: Splits the raw input CSV file into smaller, manageable chunks.
- Description:
  - ✓ Accepts a raw CSV file
  - ✓ Automatically splits into multiple files with 3 columns each (by default)
  - ✓ Helps reduce memory load and simplifies further processing
- Code Snippet:

```
def split_by_column_group(df, common_key, output_dir, split_size = None):
    for i, size in enumerate(split_sizes):
        end = start + size
        col_group = other_cols[start:end]
        sub_df = df[[common_key] + col_group]
        file_path = os.path.join(output_dir, f"split_part_{i+1}.csv")
        sub_df.to_csv(file_path, index = False)
        print(f"☑ Saved: {file_path}")
        file_paths.append(file_path)
        start = end
```

**Code Snippet 2**



### 8.3.3 Cleaning.py

- Purpose: Applies consistent cleaning rules to each dataset.
- Description:
  - ✓ Removes unwanted characters
  - ✓ Fills or drops missing values
  - ✓ Standardizes column names
  - ✓ Prepares data for further processing or storage
- Code Snippet:

```
def clean(df):  
    df.columns = df.columns.str.strip().str.lower().str.replace(' ', '_')  
    df = df.dropna(how='all') # Drop empty rows  
    df = df.fillna('N/A')    # Replace NaNs with placeholder  
    return df
```

**Code Snippet 3**



### 8.3.4 datatype\_utils.py

- Purpose: Converts pandas datatypes into MySQL-compatible formats.
- Description:
  - ✓ Automatically maps each column's dtype
  - ✓ Ensures schema compatibility for dynamic SQL generation
  - ✓ Prevents datatype mismatch during table creation
- Code Snippet:

```
def pd.api.types.is_numeric_dtype(series):
    if series.dropna().apply(lambda x:float(x).is_integer()).all():
        corrected_types[col] = 'INT'
        print("Changed Type: INT")
    else:
        corrected_types[col] = 'FLOAT'
        print("Changed Type: FLOAT")

    elif parsed.notna().mean() > 0.6:
        print(f"'{col}' is likely a date column.")
        corrected_types[col] = 'date'
    else:
        corrected_types[col] = 'VARCHAR(255)'
        print("Changed Type: VARCHAR(255)")
```

Code Snippet 4



### 8.3.5 dateconversion.py

- Purpose: Converts date values to MySQL-friendly **YYYY-MM-DD** format.
- Description:
  - ✓ Identifies date columns
  - ✓ Uses pd.to\_datetime() for safe parsing
  - ✓ Converts date to string format for MySQL compatibility
- Code Snippet:

```
try:  
    df[col] = pd.to_datetime(df[col], errors = 'coerce')  
    df[col] = df[col].dt.strftime('%Y-%m-%d')  
    print(f"{col} converted to MySQL compatible date-time format")  
except:  
    print("Failed to convert")
```

**Code Snippet 5**



### 8.3.6 read\_dfs.py

- Purpose: Loads all split CSV files into pandas DataFrames.
- Description:
  - ✓ Iterates through file list or directory
  - ✓ Reads files into a list of DataFrames
  - ✓ Prepares for batch cleaning and transformation
- Code Snippet:

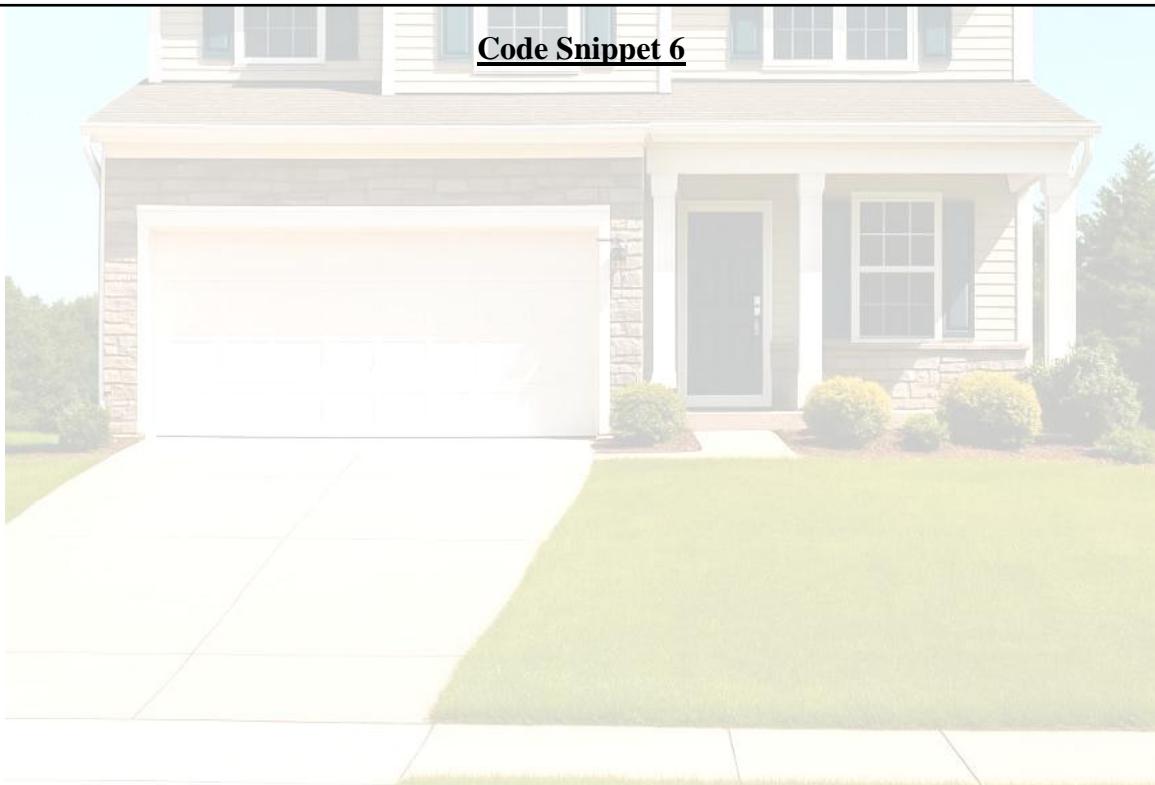
```
import pandas as pd

def read_split_files(file_paths):
    dfs = []

    for i, path in enumerate(file_paths):
        df = pd.read_csv(path)
        dfs.append(df)
        print(f"📄 Loaded {path} as df{i+1}")

    return dfs
```

**Code Snippet 6**



### 8.3.7 database\_connection.py

- Purpose: Establishes a secure connection to the MySQL database.
- Description:
  - ✓ Uses mysql.connector for connectivity
  - ✓ Handles credentials, error catching, and reconnection logic
  - ✓ Central utility for any DB I/O operation
- Code Snippet:

```
import mysql.connector

def get_connection():
    return mysql.connector.connect(
        host="localhost",
        user="your_user",
        password="your_password",
        database="real_estate_db"
    )
```

**Code Snippet 7**



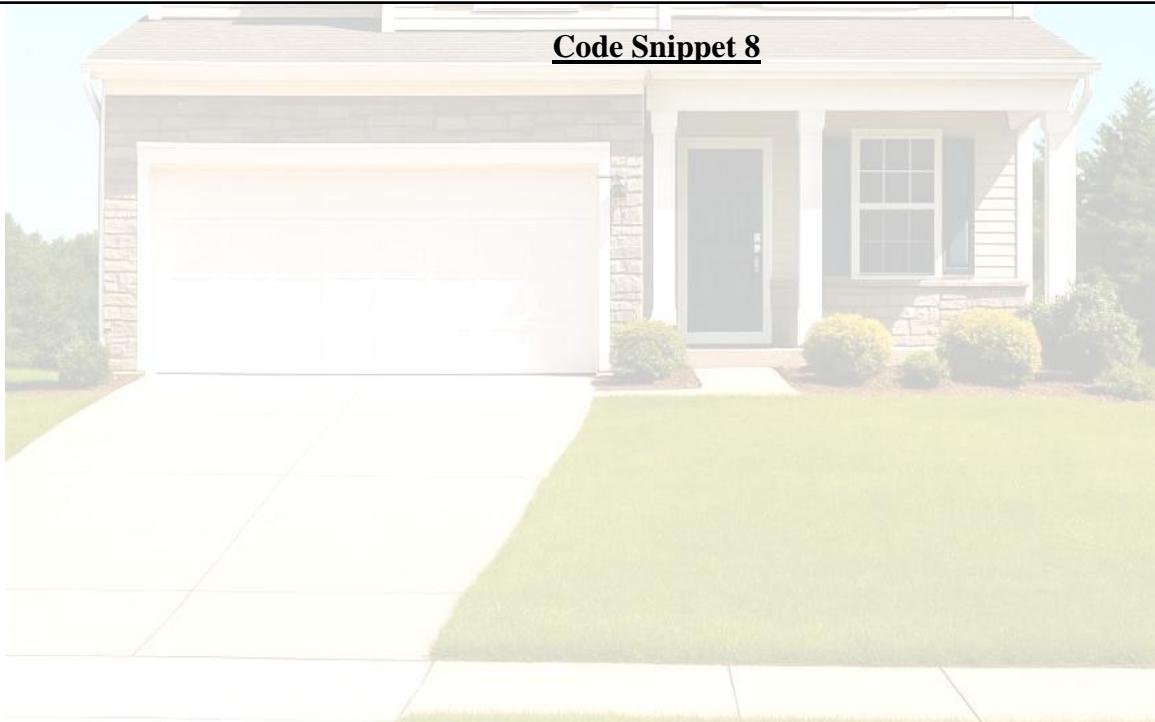
### 8.3.8 sql\_generator.py

- Purpose: Generates MySQL-compatible CREATE TABLE and INSERT INTO statements dynamically
- Description:
  - ✓ Accepts cleaned DataFrames
  - ✓ Detects column names and types
  - ✓ Creates tables and populates them with data
  - ✓ Supports dynamic schema creation
- Code Snippet:

```
# Start query
if if_not_exists:
    sql = f"CREATE TABLE IF NOT EXISTS `{table_name}` (\n"
else:
    sql = f"CREATE TABLE `{table_name}` (\n"

# Get column names and placeholders
cols = ",".join([f"`{col}`" for col in df.columns])
col_clean = cols.strip().replace(" ", "_").lower()
placeholders = ",".join(["%s"] * len(df.columns))
insert_query = f"insert into `{table_name}` ({col_clean}) VALUES ({placeholders})"
```

**Code Snippet 8**



## 9. Conclusion & Future Scope

The current ETL pipeline provides a modular, scalable, and Pythonic solution for ingesting, cleaning, transforming, and loading structured data into a MySQL database. Each utility module handles a specific responsibility, ensuring clear separation of concerns and easier maintenance.

### 9.1 Achievements

- Robust modular design using reusable Python scripts
- Dynamic table and schema generation based on input data
- Clean data transformation flow, from raw CSV to structured database tables
- Compatibility with MySQL standards for datatype and date formatting

### 9.2 Scope for Future Improvements:

#### 9.2.1 Error Logging & Monitoring:

Implement centralized logging to track issues across the ETL stages for better debugging and audit trails.

#### 9.2.2 Unit Testing:

Add unit tests for each module to ensure data integrity and reduce risk of runtime errors.

#### 9.2.3 Configuration File Integration:

Use a config file (e.g., YAML or JSON) to manage file paths, DB credentials, and schema mappings, reducing hard-coded values.

#### 9.2.4 Support for Other Data Sources:

Extend input support for JSON, Excel, or API-based data ingestion.

#### 9.2.5 Parallel Processing:

Introduce multiprocessing to speed up large dataset transformations and insertions.

#### **9.2.6 Integration with Orchestration Tools:**

Incorporate tools like Apache Airflow or Prefect for task scheduling, dependency management, and pipeline automation.

#### **9.2.7 Data Validation Layer:**

Build a schema validation check (using *pandera*, *pydantic*, or *cerberus*) before loading data into the database.



## 10. Uncovering Trends in the Real Estate Market Using Power BI

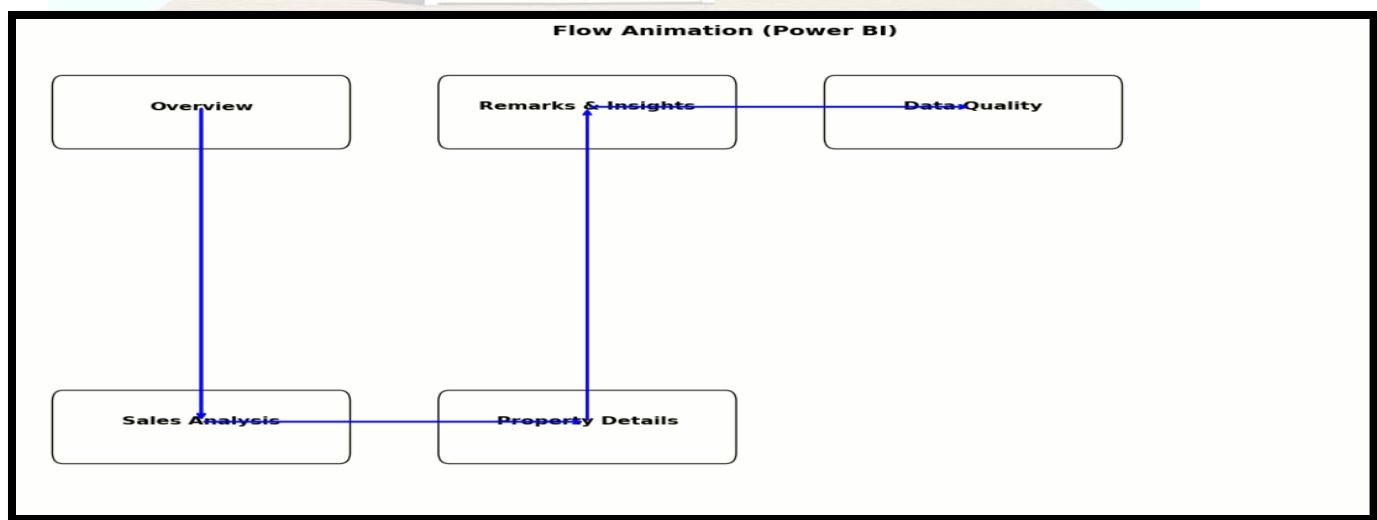
### 10.1 Introduction

The real estate industry is one of the most dynamic sectors, constantly influenced by economic shifts, buyer behavior, property characteristics, and geographic factors. This Power BI dashboard project — titled "**Uncovering Trends in the Real Estate Market**" — is designed to extract meaningful insights from historical real estate transaction data using modern data visualization techniques.

Through a series of interactive, multi-page dashboards and drillthrough reports, this project aims to provide a comprehensive view of property sales patterns, pricing dynamics, and location-based performance metrics. The goal is to empower stakeholders — from analysts to decision-makers — to explore trends, identify high-performing areas, and assess data quality with ease and clarity.

### 10.2 Dashboard Structure

The Power BI report consists of **five main dashboard pages** and **two drillthrough pages**, each designed to serve a specific analytical purpose.



**Figure 2:** Illustrates the navigation flow between dashboard sections, linking insights like Overview, Sales Analysis, and property details through interactive page transitions.

## 10.3 Main Dashboard Pages

### 10.3.1 Overview

Provides a high-level snapshot of the entire real estate dataset. This page displays key metrics such as:

- Total number of property transactions
- Total and average sale price
- Distribution of properties
- Top-performing towns based on average sale price

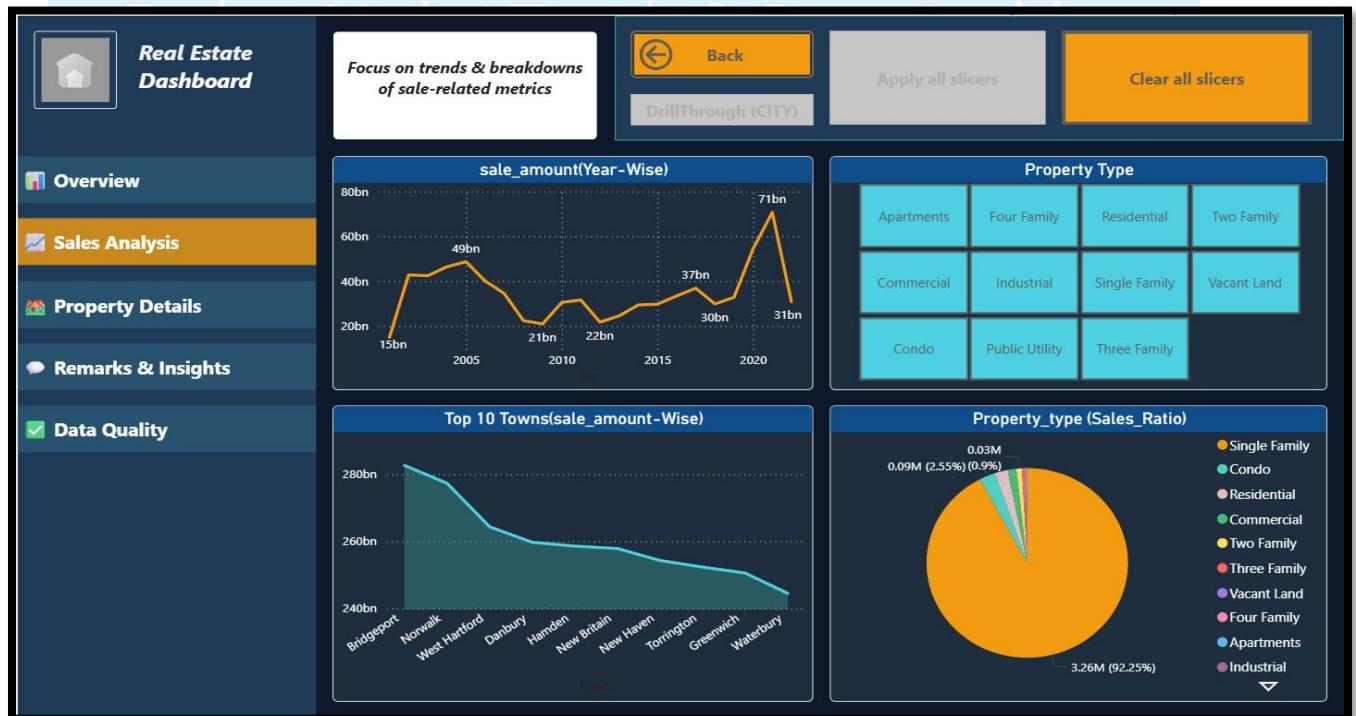


**Figure 3:** Summary of real estate data, showcasing key metrics like total property transactions (90K), average sale amount by property type, and distribution of residential properties through bar & pie charts for quick insights.

### 10.3.2 Sales Analysis

This page dives deep into sales performance metrics with:

- Monthly or yearly sales trend visualizations
- Year-wise comparisons of property sales
- Filters for town and property type
- **Drill through Link:** Connects to *Property Town Stats* for detailed location-based breakdowns



**Figure 4:** This page presents sales performance trends through year-wise line charts, top towns by sales amount and breakdown of property types, offering detailed insights into real estate sales patterns across time & locations.

### 10.3.3 Property Details

A detailed view that allows users to explore each property transaction in a tabular format.

Key features include:

- Searchable, filterable table of property listings
- Data fields such as sale date, price, address, property type, etc.
- Dynamic filters for focused exploration

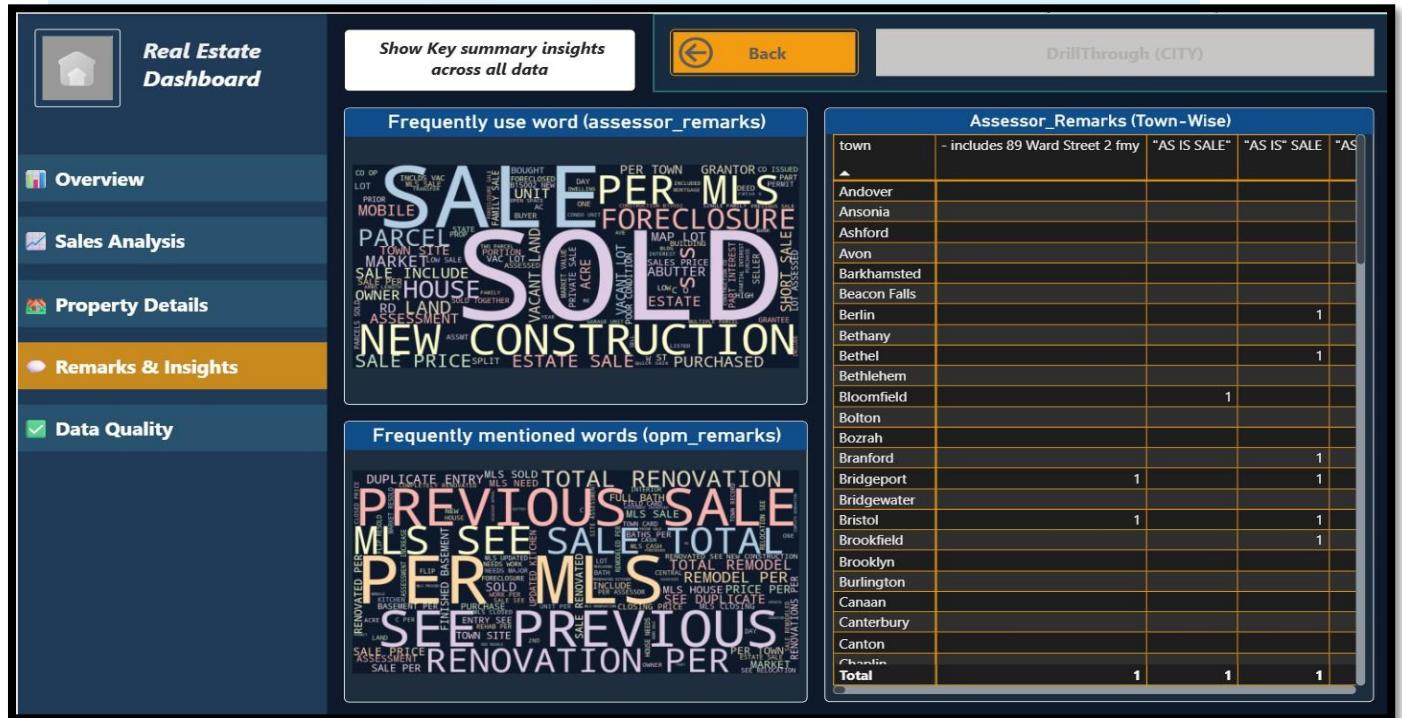


**Figure 5:** This page provides a detailed, tabular and visual breakdown of property transactions using filters, tree maps, and value tables – allowing users to analyze data by property type, usage code, sale value, and town-level distribution.

#### **10.3.4 Remarks & Insights**

This page summarizes analytical observations drawn from the data, including:

- Highlighted pricing anomalies or standout towns
  - Visual flags on inconsistent data entries
  - Insights into sale patterns, seasonality, or location-based variation
  - **Drillthrough Link:** Connects to *Property Town Stats* for deep-dive comparisons



**Figure 6:** This section highlights key data observations using word clouds and tables-surfacing pricing anomalies, frequently used remarks, and location-based trends while also flagging data inconsistencies and enabling drill-through analysis by town.

### 10.3.5 Data Quality

Dedicated to assessing the cleanliness and completeness of the dataset. This page includes:

- Null value trends
- Missing or malformed records
- Property class or zip code inconsistencies
- **Drillthrough Link:** Connects to *Residential Property Stats* for further validation and analysis



**Figure 7:** This section evaluates data completeness and accuracy using bar charts and scatter plot that display missing values across tables and fields. It highlights null trends, malformed records, and inconsistencies in property classifications, supporting better data governance and validation.

## 10.4 Drill through Pages

Drill through pages are designed to allow users to explore detailed records or breakdowns based on filters applied in main dashboards.

### 10.4.1 Property Town Stats

A focused analysis of towns linked from both:

**REAL ESTATE**  
Explore Market Trends & Opportunities

- Sales Analysis
- Remarks & Insights pages

Key insights include:

- Average sale price by town
- Sale count per town

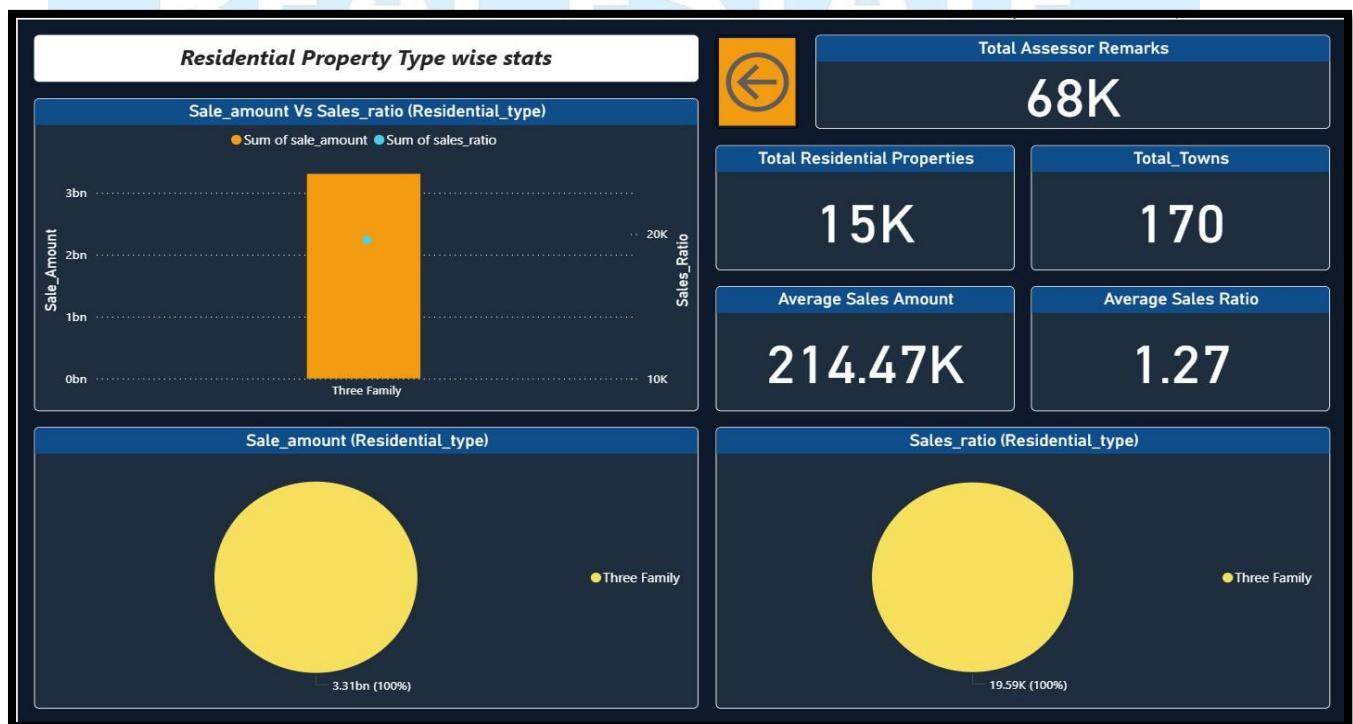


**Figure 8: Town wise Property Stats Drillthrough** – This page provides a detailed breakdown of property sales for a selected town, highlighting average sale amounts, property types, and monthly trends.

## 10.4.2 Residential Property Stats

Linked directly from the **Data Quality** page, this report provides:

- Building class-specific distribution
- Trends of residential property listings
- Comparison of sale amount and sales ratio for residential properties



**Figure 9:** This drill through page provides a detailed breakdown of residential property types, highlighting sales amount and sales ratio distributions. It supports in-depth analysis of building class trends and locality coverage across 170 towns.

## 10.5 Interactive Features

This report uses Power BI's powerful interactive features to enhance usability:

- **Cross-filtering and slicers** to dynamically adjust visuals
- **Drillthrough navigation** for focused contextual analysis
- **KPI cards** for real-time metric updates
- **Hover tooltips** for quick reference insights
- **Tab-level bookmarks or buttons** for intuitive navigation between pages

## 10.6 Value Proposition Market Trends & Opportunities

By leveraging the analytical capabilities of Power BI, this project offers:

- An exploratory data tool for understanding sales performance
- Visual storytelling for real estate stakeholders
- A benchmark to assess data quality and readiness for predictive analytics
- A modular structure that can be extended with new metrics or real-time integrations in future iterations



## 11. Discussion

- The project combines a robust **Python-based ETL pipeline**, **MySQL database integration**, and **Power BI dashboards** to explore real estate sales trends from 2001 to 2022.
- The modular ETL process (splitting, cleaning, formatting, and loading) ensures:
  - ✓ Scalability for large datasets
  - ✓ Flexibility for schema changes
  - ✓ Automation of SQL script generation and data loading
- **Power BI dashboards** provide an interactive and multi-page interface that highlights:
  - ✓ Yearly and regional property sale patterns
  - ✓ Key metrics like total sales, average amount, and property distribution
  - ✓ Drill through functionality for detailed town-wise and residential-type stats
- The **Property Town Stats** and **Residential Property Stats** pages allow users to explore:
  - ✓ Localized market activity
  - ✓ Sales distribution across building types
  - ✓ Data accuracy and completeness
- This architecture supports:
  - ✓ Business insights for real estate planners and investors
  - ✓ Data quality verification through dedicated reporting
- Overall, the solution demonstrates how **data engineering + visualization** can drive **data-informed decision-making** in the real estate domain.

## 12. Limitations

- The dataset is historical (2001–2022) and may not reflect current market dynamics post-2022.
- Property attribute coverage is limited—data like property condition, owner history, or real-time valuations are missing.
- ETL pipeline assumes consistent input format; unexpected schema changes or corrupt files may require manual intervention.
- Power BI dashboards are dependent on static snapshot data; no real-time data streaming is implemented.
- Drill through insights are limited to pre-configured fields—advanced AI-driven filters or clustering are not included.
- Currently optimized for a single database (MySQL); not tested with other RDBMS or cloud-based storage systems.

## 13. Future Scope

- Extend the ETL pipeline to **ingest live streaming data** using APIs or web scraping for real-time updates.
- Incorporate **machine learning models** to predict future property prices, identify outlier sales, or recommend investment locations.
- Improve **error handling and logging** in the ETL modules for production-grade robustness.
- Add support for **multi-database storage options**, including PostgreSQL and cloud databases like AWS RDS or Google BigQuery.
- Enhance dashboard interactivity with **AI-based Q&A** visual and **user-driven filters**.
- Introduce **user roles and access control** in Power BI Service for secure report sharing and governance.

## 14. Conclusion

- Successfully developed an **end-to-end ETL pipeline** using Python to clean, transform, and load real estate data into a structured MySQL database.
- Designed a modular architecture with reusable utility scripts, promoting **scalability, maintainability, and efficiency**.
- Built an insightful **Power BI dashboard** comprising five main pages and two drill through pages, enabling multidimensional analysis of sales, properties, and trends.
- Provided detailed **visual exploration of market behavior**, including sales distribution, property types, city-level trends, and data quality indicators.
- Demonstrated how **data integration, automation, and visualization** can deliver real business value in real estate analytics.
- Created a strong foundation for expanding the system with **advanced analytics, machine learning models, or live data feeds** in the future.

## 16. Analyst Profile

- **Name:** Anurag Yadav
- **Role:** Data Analyst | Business Intelligence Developer
- **Technical Skills:** Power BI, DAX, SQL, MySQL, Python (pandas, numpy, seaborn, streamlit), Excel, Data Cleaning, Data Modelling, Data Visualization
- **Industry Experience:** 7+ years in the IT industry with hands-on experience in data analysis, reporting, and visualization solutions.
- **Project Focus:** Skilled in transforming raw data into actionable insights through interactive dashboards and reports. Passionate about uncovering patterns in user behavior and improving decision-making through data.
- **Tools Used in This Project:** Power BI Desktop, DAX, Power Query
- **Professional Interests:** Data storytelling, user engagement analytics, dashboard design, predictive analytics, and music data analysis.

\*\*\*