# Data Collection and Preprocessing Phase

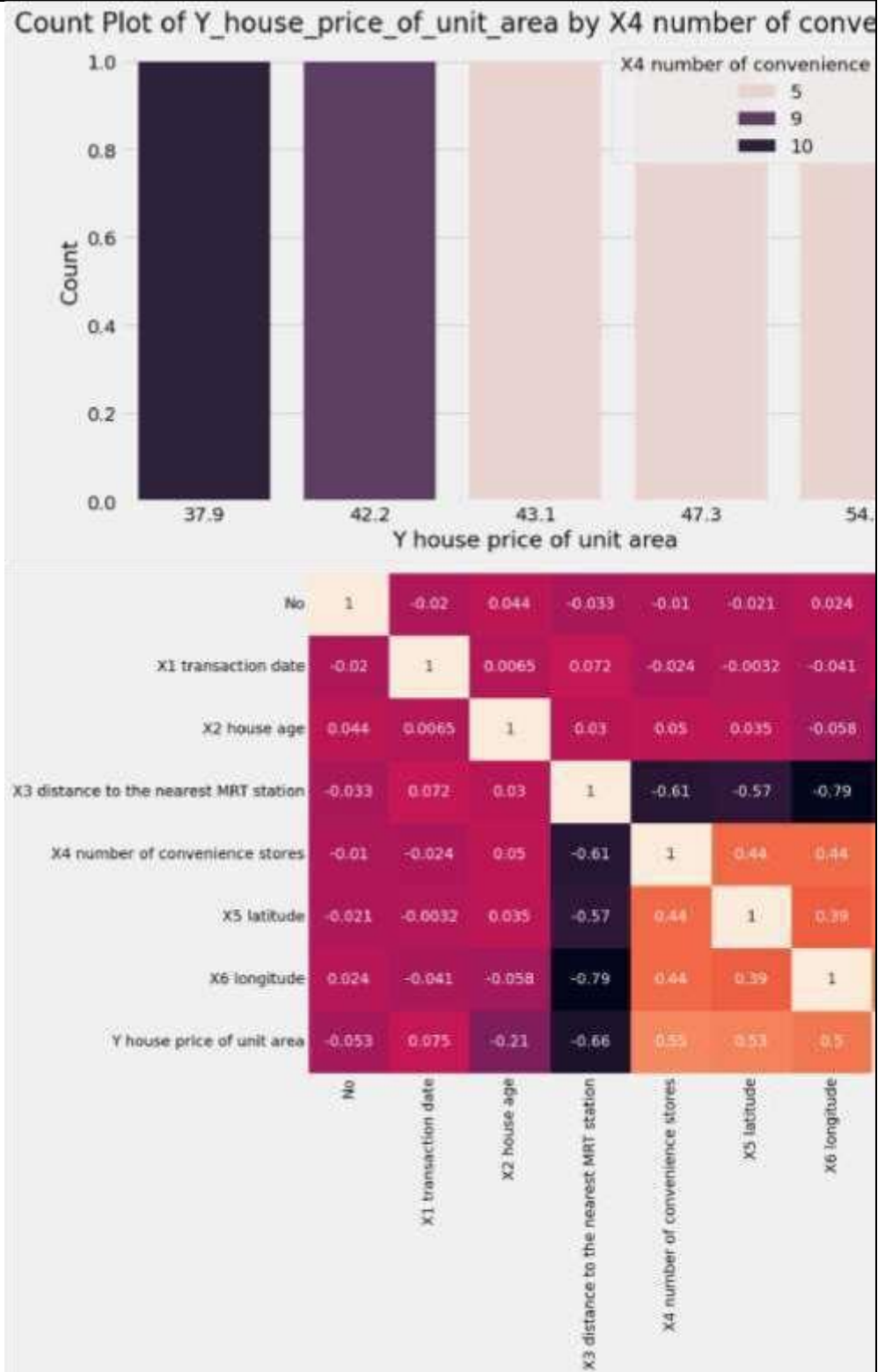| | |
|---|---|
| Date | 8 July 2024 |
| Team ID | 740138 |
| Project Title | Identification Of Methodology Used In Real Estate Property Valuation |
| Maximum Marks | 6 Marks |

**Data Exploration and Preprocessing Template**

Identifies data sources, assesses quality issues like missing values and duplicates, and implements resolution plans to ensure accurate and reliable analysis.
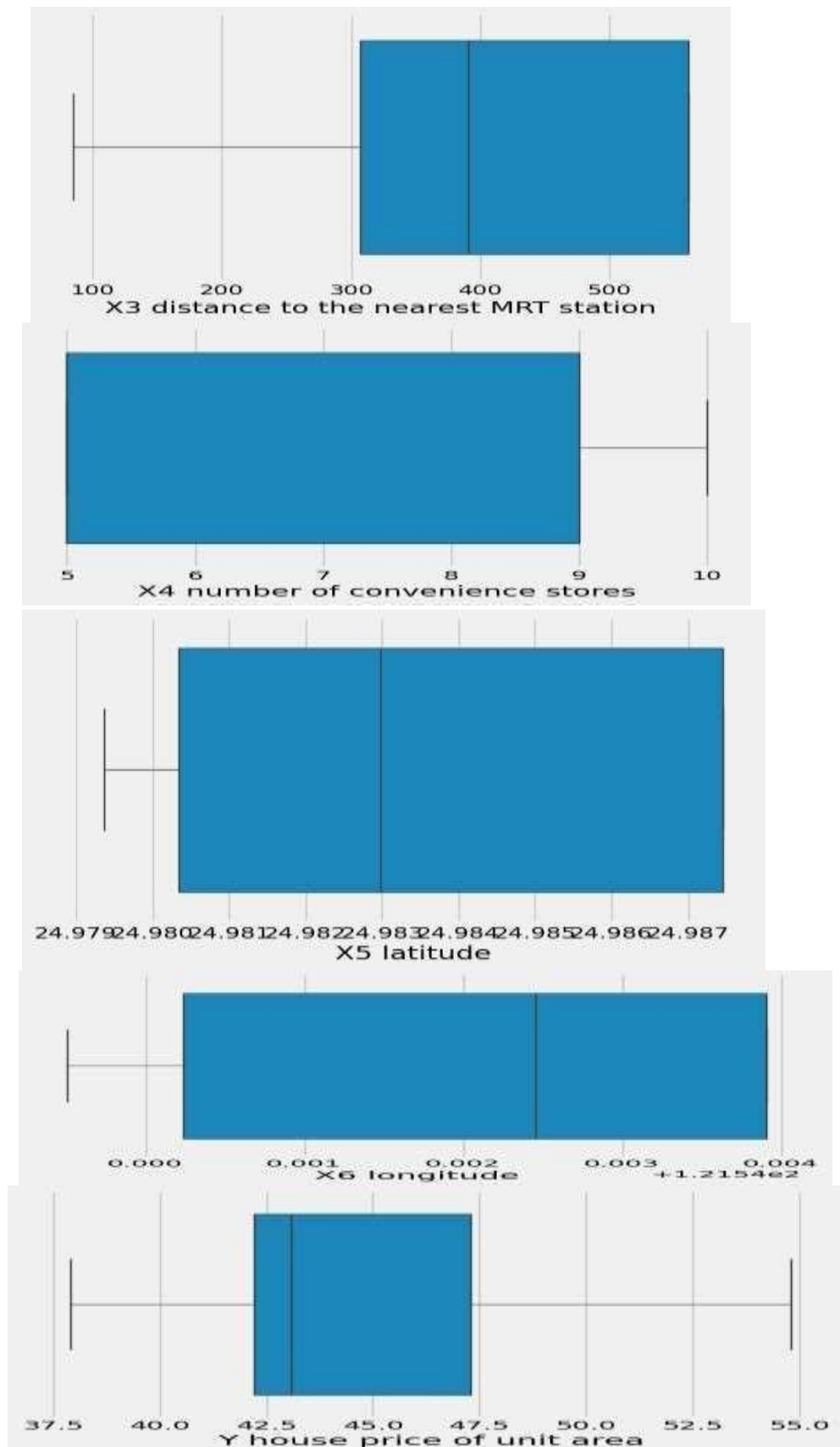
| Section | Description |
|---|---|
| Data Overview | Dimension: <br> 331rowsx 8columns Descriptive statistics: <br>  |
| Univariate Analysis |  |

| | |
|---|---|
| |  |
| Bivariate Analysis | 
Scatterplot of X3 distance to the nearest MRT station vs X4 number of conve |

| Multivariate Analysis | |
|---|---|

**Count Plot of Y_house_price_of_unit_area by X4 number of conve**

X4 number of convenience
- 5
- 9
- 10

Count

1.0
0.8
0.6
0.4
0.2
0.0

37.9   42.2   43.1   47.3   54.

Y house price of unit area

| | No | X1 transaction date | X2 house age | X3 distance to the nearest MRT station | X4 number of convenience stores | X5 latitude | X6 longitude |
|---|---|---|---|---|---|---|---|
| No | 1 | -0.02 | 0.044 | -0.033 | -0.01 | -0.021 | 0.024 |
| X1 transaction date | -0.02 | 1 | 0.0065 | 0.072 | -0.024 | -0.0032 | -0.041 |
| X2 house age | 0.044 | 0.0065 | 1 | 0.03 | 0.05 | 0.035 | -0.058 |
| X3 distance to the nearest MRT station | -0.033 | 0.072 | 0.03 | 1 | -0.61 | -0.57 | -0.79 |
| X4 number of convenience stores | -0.01 | -0.024 | 0.05 | -0.61 | 1 | 0.44 | 0.44 |
| X5 latitude | -0.021 | -0.0032 | 0.035 | -0.57 | 0.44 | 1 | 0.39 |
| X6 longitude | 0.024 | -0.041 | -0.058 | -0.79 | 0.44 | 0.39 | 1 |
| Y house price of unit area | -0.053 | 0.075 | -0.21 | -0.66 | 0.55 | 0.53 | 0.5 |

| | |
|---|---|
| Handled Outliers and Anomalies |  |

**Data Preprocessing Code Screenshots**

| | |
|---|---|
| Loading Data |  |
| Finding &Handling Missing Data |  |
| Data Transformation | - |
| Feature Engineering | Attached the code in final submission |
| Save Processed Data |  |