# Data Collection and Preprocessing Phase

| Date | 8 July 2024 |
|---|---|
| Team ID | 740109 |
| 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.
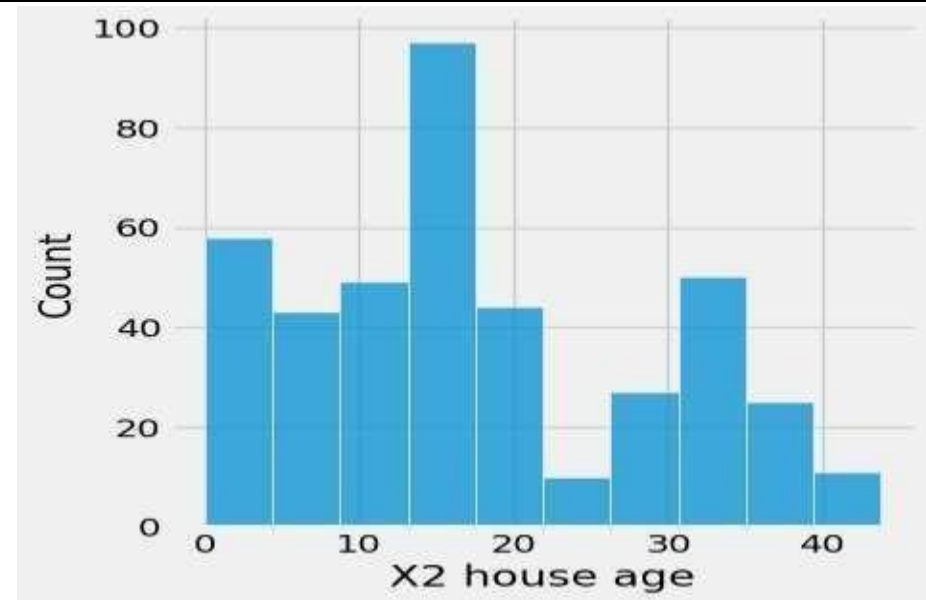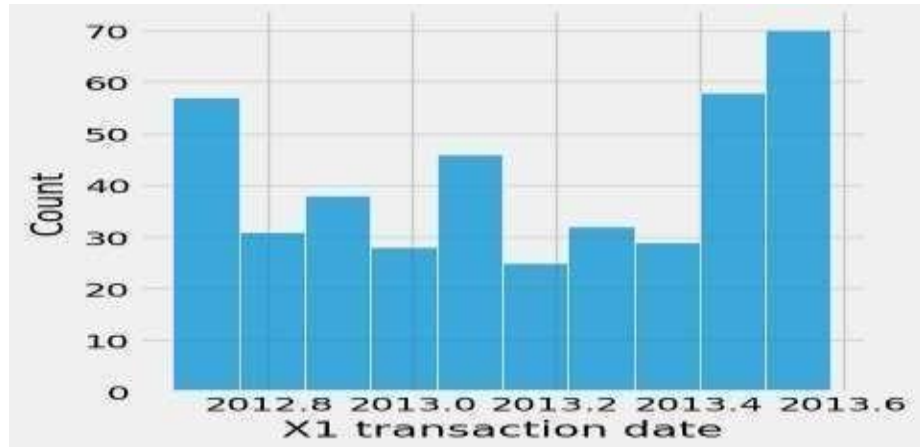
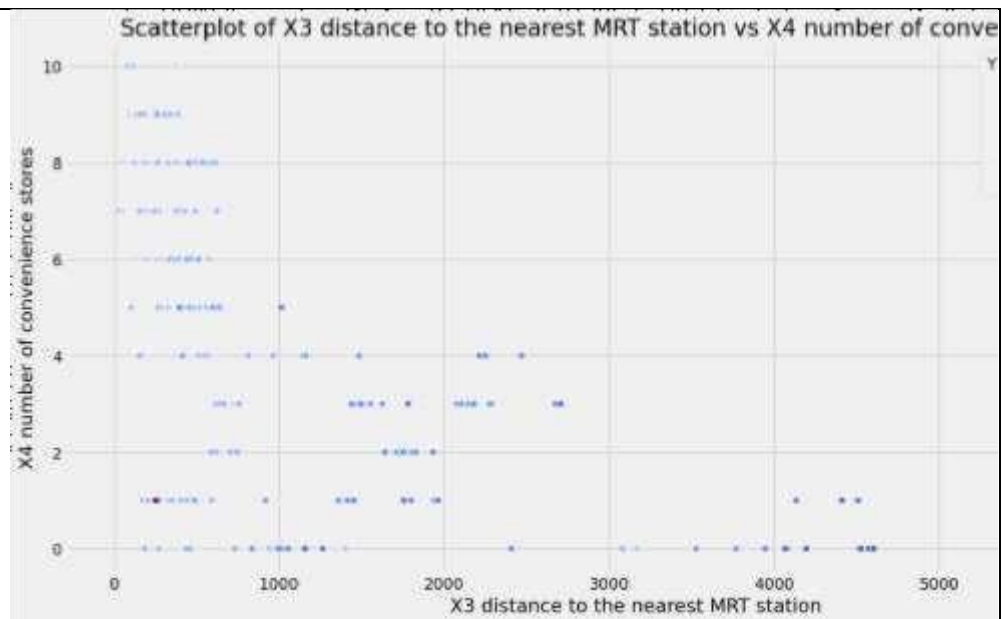| Section | Description |
|---|---|
| Data Overview | Dimension:<br>331rowsx 8columns Descriptive statistics:<br> |

| 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

| | |
|---|---|
| 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 | ```python
import pickle
from sklearn.preprocessing import StandardScaler
from sklearn.ensemble import RandomForestRegressor
rf_model = RandomForestRegressor()
scaler = StandardScaler()
with open('price.pkl', 'wb') as f:
    pickle.dump(rf_model, f)
with open('scale.pkl', 'wb') as f:
    pickle.dump(scaler, f)
```

```python
from google.colab import files
files.download('price.pkl')
```

```python
from google.colab import files
files.download('scale.pkl')
```

```python
from google.colab import files
files.download('/content/drive/MyDrive/ dataset/real estate valuation data set.csv')
``` |