```
# Mount Google Drive to access the dataset
from google.colab import drive
drive.mount('/content/drive')
# File path to the dataset
file_path = '/content/drive/MyDrive/House Price Prediction & Feature Impact Analysis/Real es
# Check if the file exists
import os
print("File exists:", os.path.exists(file_path))
# Load the dataset into a Pandas DataFrame
import pandas as pd
data = pd.read_csv(file_path)
# Display the first few rows and dataset shape
print(data.head())
print("Dataset shape:", data.shape)
→▼ Mounted at /content/drive
     File exists: True
        No X1 transaction date X2 house age \
     0
         1
                       2012.917
                                          32.0
     1
         2
                       2012.917
                                          19.5
     2
         3
                       2013.583
                                          13.3
         4
     3
                       2013.500
                                          13.3
         5
                       2012.833
                                           5.0
        X3 distance to the nearest MRT station X4 number of convenience stores
     0
                                       84.87882
                                                                               10
     1
                                      306.59470
                                                                                9
     2
                                                                                5
                                      561.98450
     3
                                                                                5
                                      561.98450
     4
                                                                                5
                                      390.56840
        X5 latitude X6 longitude Y house price of unit area
     0
           24.98298
                        121.54024
                                                          37.9
     1
           24.98034
                        121.53951
                                                          42.2
                                                          47.3
     2
           24.98746
                        121.54391
           24.98746
     3
                        121.54391
                                                          54.8
           24.97937
                        121.54245
                                                          43.1
     Dataset shape: (414, 8)
# Rename columns for better readability
data.columns = [
    "ID", "Transaction_Date", "House_Age",
    "Distance_to_MRT", "Convenience_Stores",
    "Latitude", "Longitude", "Price_per_Unit_Area"
]
```

```
# Convert Transaction_Date to a datetime format and extract year
data['Transaction_Date'] = pd.to_datetime(data['Transaction_Date'], format='%Y.%f', errors='
data['Transaction_Year'] = data['Transaction_Date'].dt.year
# Remove outliers from the Price_per_Unit_Area column using the IQR method
def remove outliers(df, column):
    Q1 = df[column].quantile(0.25)
    Q3 = df[column].quantile(0.75)
    IQR = Q3 - Q1
    lower_bound = Q1 - 1.5 * IQR
    upper bound = Q3 + 1.5 * IQR
    return df[(df[column] >= lower_bound) & (df[column] <= upper_bound)]</pre>
data = remove_outliers(data, 'Price_per_Unit_Area')
# Create distance categories
import numpy as np
data['Distance_Category'] = pd.cut(
    data['Distance_to_MRT'],
    bins=[0, 500, 1000, 3000, 5000, np.inf],
    labels=['Very Close', 'Close', 'Moderate', 'Far', 'Very Far']
)
# Save the cleaned dataset
processed_file_path = '/content/drive/MyDrive/House Price Prediction & Feature Impact Analys
data.to csv(processed file path, index=False)
print("Processed data saved at:", processed_file_path)
→ Processed data saved at: /content/drive/MyDrive/House Price Prediction & Feature Impact
     <ipython-input-2-833b4a2ec6b2>:25: SettingWithCopyWarning:
     A value is trying to be set on a copy of a slice from a DataFrame.
     Try using .loc[row_indexer,col_indexer] = value instead
     See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user
       data['Distance_Category'] = pd.cut(
import seaborn as sns
import matplotlib.pyplot as plt
# Plot house price distribution
plt.figure(figsize=(8, 6))
sns.histplot(data['Price_per_Unit_Area'], kde=True, bins=30)
plt.title('Distribution of House Prices')
plt.xlabel('Price per Unit Area')
plt.ylabel('Frequency')
plt.show()
# Plot feature correlations
plt.figure(figsize=(10, 8))
```

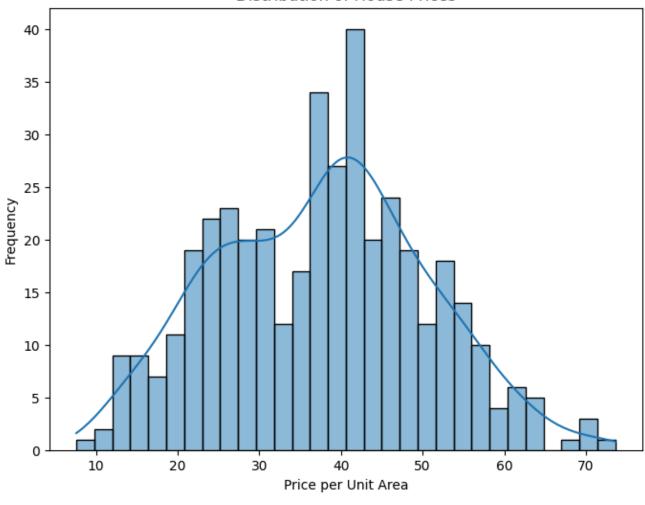
```
sns.heatmap(data.select_dtypes(include=['float64', 'int64']).corr(), annot=True, cmap='coolw
plt.title('Feature Correlation Heatmap')
plt.show()

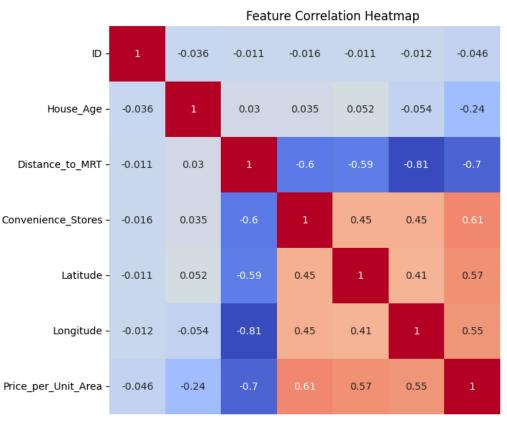
# Scatter plots for key relationships
plt.figure(figsize=(8, 6))
sns.scatterplot(x='Distance_to_MRT', y='Price_per_Unit_Area', data=data)
plt.title('Price vs. Distance to MRT')
plt.show()

plt.figure(figsize=(8, 6))
sns.scatterplot(x='House_Age', y='Price_per_Unit_Area', data=data)
plt.title('Price vs. House Age')
plt.show()
```



Distribution of House Prices





1.00

Transaction_Year -

