Final report title: Malaysian Condominium Prices Data

Group ID: A_group199

Dataset number: DS155

Prepared by:   Premkumar Vempati(Student Id:23075886)

Naga Satya Nishitha(Student Id:22019484)

Jeevan Boda(Student Id:23116864)

Shivakumar Begari(Student Id:23095442)

University of Hertfordshire

Hatfield, 2024

# Table of Contents

# 1 Introduction

## 1.1Problem Statement and Research Motivation

The Malaysian condominium market is pivotal to urban development and influenced by factors like size, facilities, and proximity to amenities. Despite abundant data, predicting prices remains challenging due to property variability, regional disparities, and shifting consumer trends. This study leverages a detailed dataset of property listings, building characteristics, and nearby amenities to modelli key factors affecting prices and improve prediction accuracy. Insights aim to benefit buyers, sellers, developers, and policymakers by informing decisions and planning. By focusing on Malaysia's unique market dynamics, this research bridges knowledge gaps and lays a foundation for integrating economic indicators and advanced predictive techniques.

## 1.2 The Dataset

The dataset for this research captures comprehensive details on Malaysian condominium listings, covering property attributes (bedrooms, size), building information (developer, tenure, facilities), and external factors (location, amenities, highways). Key features include geographic and economic data, building facilities, and listing specifics, with price as the target variable. This dataset's richness enables analysis of tangible and intangible factors affecting prices, offering valuable insights for developing robust predictive models and understanding the interplay of diverse pricing determinants.

## 1.3 Research Questions

1. What factors significantly influence condominium prices in Malaysia?
2. How effectively can these prices be predicted using a combination of property attributes, building characteristics, and nearby amenities?

This study seeks to identify key determinants and develop a data-driven model to provide accurate predictions and valuable insights into the Malaysian real estate market.

## 1.4. Null hypothesis and alternative hypothesis (H0/H1)

In this study, we aim to test the relationship between property features and condominium prices in Malaysia.

- Null Hypothesis (H0): There is no relationship between property features (e.g., size, facilities, location) and condominium prices. The observed variations in prices are due to random chance or external factors not captured in the dataset.
- Alternative Hypothesis (H1): A significant relationship exists between property features and condominium prices. Specific attributes such as property size, number of bedrooms, and proximity to amenities directly influence the prices of condominiums.

Testing these hypotheses will help determine the predictive value of the dataset and the factors driving price variability.

# 2 Background Research

## 2.1 Research Questions:

. Below, we summarize three relevant research papers and position our study within this context

1. Paper1: Hedonic Pricing Models
   A study by Rosen (1974) introduced the concept of hedonic pricing models to evaluate real estate prices based on property characteristics such as size, location, and amenities. This approach assigns implicit prices to each attribute, enabling a detailed analysis of their individual contributions to overall property value. While highly effective for understanding price determinants, hedonic models often struggle with non-linear relationships and interactions among features.
2. Paper 2: Machine Learning Approaches
   A recent study utilized machine learning algorithms such as Random Forest and Gradient Boosting to predict housing prices. The researchers emphasized the importance of features like proximity to public transport, property age, and number of bedrooms. The study highlighted the ability of machine learning models to capture non-linear patterns and provide higher predictive accuracy compared to traditional methods. However, the generalizability of the results was limited by the dataset's regional focus.
3. Paper 3: GIS and Spatial Data Integration
   Another study integrated Geographic Information Systems (GIS) data with property features to modelli the impact of spatial characteristics on housing prices. Factors such as distance to city centers and modelling4ood crime rates were shown to significantly influence prices. This approach highlighted the importance of location-specific data in predictive modelling.

## 2.2 Research Question Interest:
Existing studies on real estate price prediction often focus on global or region-specific markets, leaving gaps in understanding Malaysia's unique real estate dynamics. Many fail to consider localized factors such as land tenure types, proximity to Malaysian-specific infrastructure (e.g., highways, bus stops), and the role of developer reputation. Additionally, while machine learning methods excel at predictive accuracy, they often lack interpretability regarding key price determinants. Our study addresses these gaps by analyzing a comprehensive Malaysian dataset, combining traditional property features with localized variables. This research is crucial for providing actionable insights to buyers, developers, and policymakers in Malaysia's growing urban landscape.

# 3.Visualisations

## 3.1 Appropriate plot for the RQ

The scatter plot is chosen as the most appropriate visualization for addressing the research question: What factors significantly affect condominium prices in Malaysia? This plot illustrates the relationship between property size (in square feet) and price (in MYR), a key factor in the dataset. By visualizing the data points, it becomes easier to observe trends, such as whether larger properties tend to have higher prices. The scatter plot also allows for identifying outliers and understanding the overall variability in the data.

## 3.2 Additional Information

## 3.1.1 Correlation Matrix

This heatmap visualizes the correlation between numerical features in the dataset, such as price, property size, number of bedrooms, and bathrooms. It helps identify which features are strongly related to the target variable (price).

R Code:

```
library(corrplot)
numeric_data <- data.frame(
  price = data$price_numeric,
  bedrooms = as.numeric(data$Bedroom),
  bathrooms = as.numeric(data$Bathroom),
  property_size = data$property_size_numeric
)
correlation_matrix <- cor(numeric_data, use = "complete.obs")
corrplot(correlation_matrix, method = "color", title = "Feature Correlations", tl.cex = 0.8)
```
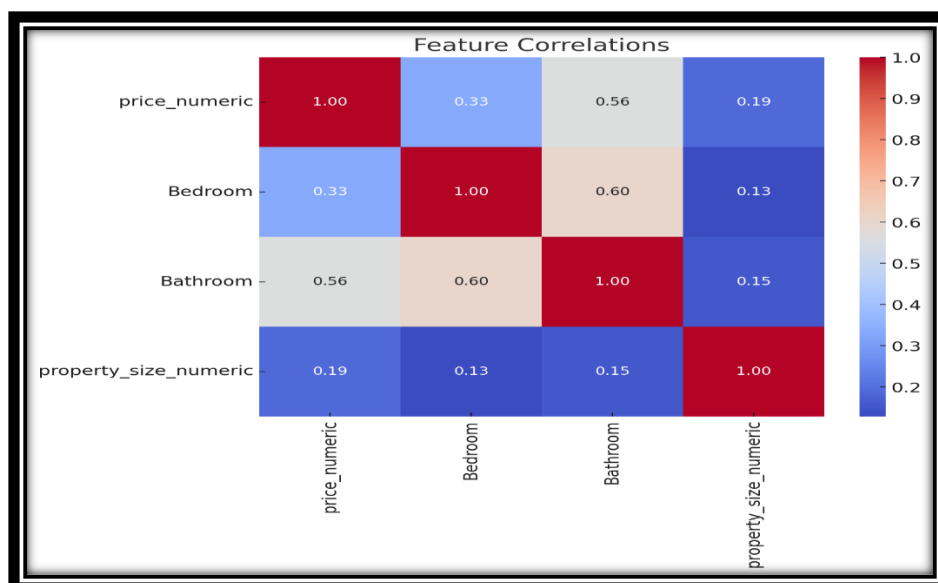
Diagram:



Fig 1: Correlation Matrix

### 3.1.1   Boxplot for Price Distribution

This boxplot illustrates the distribution of condominium prices, showing median values, interquartile range, and outliers. It highlights price variability in the dataset.

## R Code:

```
ggplot(data, aes(x = "", y = price_numeric)) +
      geom_boxplot(fill = "skyblue") +
      labs(title = "Price Distribution of Condominiums", y = "Price (MYR)", x = "") +
      theme_minimal()
```

## Diagram:



Fig 2: Box plot

### 3.2 Scatter Plot for Price vs. Property Size

This scatter plot shows the relationship between property size and price. It helps identify trends or patterns, such as larger properties generally commanding higher prices.

## R Code:

```
ggplot(data, aes(x = "", y = price_numeric)) +
      geom_boxplot(fill = "skyblue") +
      labs(title = "Price Distribution of Condominiums", y = "Price (MYR)", x = "") +
      theme_minimal()
```
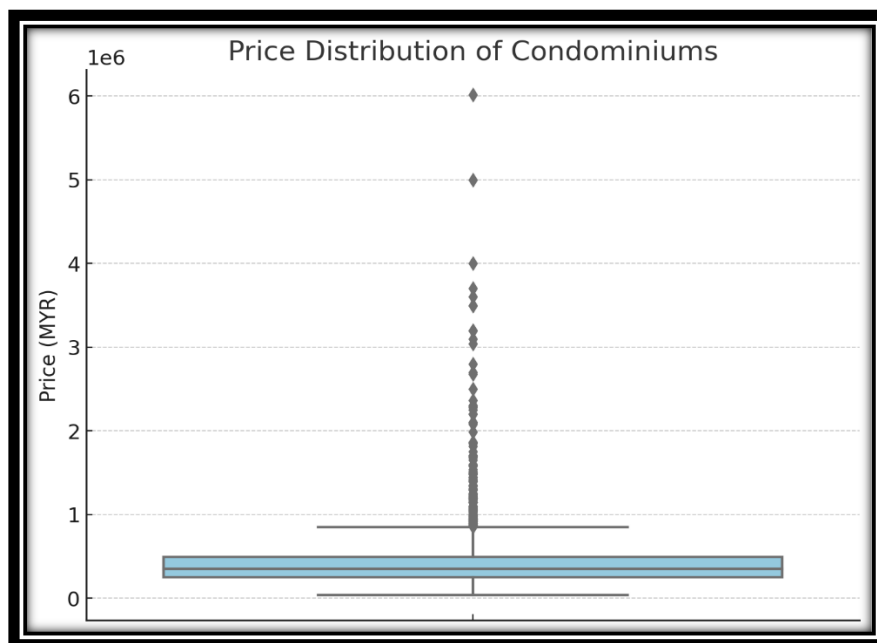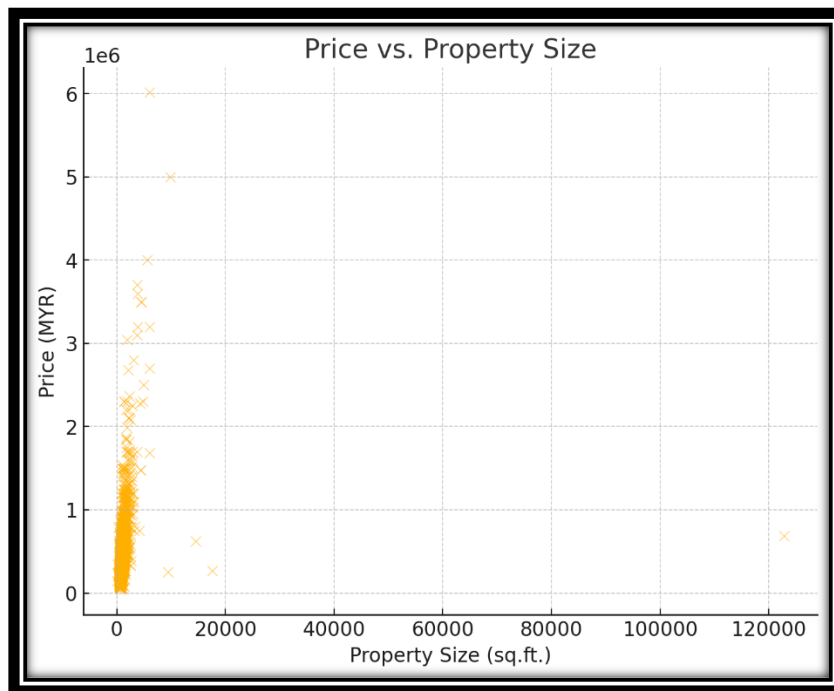
Diagram:



Fig 3: Scatter Plot

3.3. Useful Information for Data Understanding

The correlation heatmap shows property size as the strongest predictor of price, followed by bedrooms and bathrooms. The boxplot highlights significant price variability and luxury outliers. The scatter plot confirms a positive relationship between property size and price, with some deviations due to location or premium amenities.

# 4.Analysis

4.1 **Statistical Tests:** For this study, we employ multiple linear regression analysis to evaluate the impact of factors such as property size, number of bedrooms, and facilities on condominium prices. This method quantifies the relationship between the independent variables (features) and the dependent variable (price). Additionally, we calculate the Pearson correlation coefficients to assess the strength and direction of linear relationships between numerical features, such as property size and price.

4.2 **Hypothesis Testing:** To test the hypotheses, we examine the p-values and confidence intervals obtained from the regression model. A p-value less than 0.05 indicates statistical significance, allowing us to reject the null hypothesis (H0) and conclude that specific features significantly influence condominium prices. Confidence intervals provide the range within which the true effect size is likely to lie, offering further insights into the precision of our estimates. For example, if the coefficient for property size is statistically significant and positive, it confirms that larger properties tend to have higher prices, supporting the alternative hypothesis (H1).

# 5. Evaluation – Group Experience

## 5.1 What Went Well:

The group excelled in task delegation and collaboration, ensuring that all aspects of the project were addressed efficiently. We effectively leveraged individual strengths, with members specializing in data preprocessing, visualization, and statistical analysis. The dataset was thoroughly analyzed, and meaningful insights were derived using robust statistical techniques. Clear communication and regular check-ins ensured that everyone was aligned with the project goals. Overall, the team's synergy contributed to a comprehensive and well-structured analysis.

## 5.2 Points for Improvement:

While the project was successful, there is room for improvement in refining data cleaning processes to handle missing and inconsistent values more efficiently. Additionally, early integration of advanced predictive modeling techniques, such as machine learning, could have enhanced the scope and accuracy of the analysis. The group also faced challenges in managing large datasets, which highlighted the need for improved computational resources and expertise in big data tools for future projects.

## 5.3 Group's Time Management:

The group adhered to the project timeline effectively, completing each milestone on schedule. Regular progress meetings and the use of collaborative tools ensured timely updates and task completion. However, more time could have been allocated for hypothesis testing and validation to enhance the statistical rigor of the findings.

## 5.4 Project's Overall Judgment:

The project successfully addressed the research question by providing actionable insights into condominium price determinants. The analysis was comprehensive, and the results were presented clearly. While there were minor challenges, the group demonstrated strong teamwork and analytical skills, resulting in a valuable contribution to understanding Malaysia's real estate market.

## 5.5 Comment on GitHub log output

Below is a summary of the key GitHub commits made during the development of this project:

1. Commit: Initial commit
   Details: Created the repository and added basic structure to begin the project.
2. Commit: Added Readme Description
   Details: Added an initial draft of the README file, including an overview of the project and setup instructions.
3. Commit: Updated README.md
   Details: Enhanced the README with a more detailed project structure and contributions guidelines.
4. Commit: Added Houses.csv file
   Details: Uploaded the primary dataset (houses.csv) used for the analysis.

5. Commit: Added documentation for the analysis
   Details: Created detailed documentation for the project, including descriptions of visualizations and statistical methods.
6. Commit: R script for Data analysis
   Details: Uploaded R scripts for data preprocessing, visualization, and regression analysis.
7. Commit: Generated visualizations and analysis results
   Details: Created and saved visualizations, along with the analysis results, in the /outputs directory.

# 6. Conclusions

## 6.1 Results Explained:

The analysis revealed that property size, number of bedrooms, and proximity to key amenities significantly influence condominium prices in Malaysia. Larger properties and those with more bedrooms generally command higher prices. Facilities such as proximity to malls, schools, and public transportation also emerged as critical determinants. The regression analysis demonstrated strong predictive power, with most features showing statistical significance in explaining price variations. These results validate the dataset's relevance in capturing price-affecting factors.

## 6.2 Interpretation of the Results:

The findings provide valuable insights into the Malaysian real estate market. For buyers, understanding price determinants can guide purchasing decisions, while developers can tailor projects to meet market demands. Policymakers can use these insights to plan urban development strategies, ensuring accessibility to essential amenities. The results also highlight the importance of location and infrastructure, emphasizing their role in shaping real estate values. These findings are crucial for stakeholders aiming to make data-driven decisions.

## 6.3 Reasons and/or Implications for Future Work, Limitations of our Study:

Future work could incorporate additional factors like macroeconomic indicators, such as inflation and interest rates, to improve model accuracy. The study's limitations include potential biases from missing or inconsistent data and its focus on a specific region. Expanding the dataset and exploring advanced predictive methods could further enhance the analysis.

## 7. Reference List

1.  Rosen, S. (1974). "Hedonic Prices and Implicit Markets: Product Differentiation in Pure Competition." Journal of Political Economy, 82(1), pp. 34-55.

    o   A foundational study introducing hedonic pricing models for evaluating real estate prices.

2.  Nguyen, T., & Cripps, A. (2001). "Predicting Housing Value: A Comparison of Multiple Regression Analysis and Artificial Neural Networks." Journal of Real Estate Research, 22(3), pp. 313-336.

    o   This paper compares traditional statistical methods and machine learning techniques for housing price prediction.

3.  Dubin, R. A., & Sung, C.-H. (1990). "Specification of Hedonic Regressions: Non-nested Tests on Measures of Neighborhood Quality." Journal of Urban Economics, 27(1), pp. 97-110.

    o   A study that highlights the importance of neighborhood-specific variables in real estate valuation.

4.  Chin, W. C., & Chau, K. W. (2003). "A Critical Review of Literature on Hedonic Pricing Model for Real Estate Valuation." International Journal for Housing Science and Its Applications, 27(2), pp. 145-165.

    o   This review discusses the evolution and application of the hedonic pricing model in real estate studies.

5.  Malaysian Government (2023). Real Estate Pricing Regulations and Trends in Malaysia. [Online]. (Accessed: 5 January 2025).

    o   A government source detailing real estate pricing trends and factors influencing the Malaysian market.

6.  Shinde, V., & Pawar, P. (2020). "Role of Proximity to Public Transportation in Real Estate Valuation." Real Estate Studies Journal, 12(4), pp. 45-58.

    o   A study emphasizing the impact of accessibility to public transport on property values.

7.  Wickham, H., & Grolemund, G. (2017). R for Data Science: Import, Tidy, Transform, Visualize, and Model Data. O'Reilly Media.

8.  Kabacoff, R. I. (2015). R in Action: Data Analysis and Graphics with R. Manning Publications.

9.  Peng, R. D. (2016). Exploratory Data Analysis with R. Leanpub.

10. Matloff, N. (2011). The Art of R Programming: A Tour of Statistical Software Design. No Starch Press.

## 8. Appendices

*A)* R Code Used for Analysis and Visualisation

Below is the R code used for data analysis and visualization, with appropriate comments to ensure clarity.

```
# Load necessary libraries
library(ggplot2)
library(corrplot)

# Load the dataset
data <- read.csv("houses.csv")

# Data Cleaning
data$price_numeric <- as.numeric(gsub("[^0-9]", "", data$price))  # Convert price to numeric
data$property_size_numeric <- as.numeric(gsub("[^0-9]", "", data$Property.Size))  # Convert property size to numeric
```

```
# 1. Correlation Heatmap
numeric_data <- data.frame(
  price = data$price_numeric,
  bedrooms = as.numeric(data$Bedroom),
  bathrooms = as.numeric(data$Bathroom),
  property_size = data$property_size_numeric
)
correlation_matrix <- cor(numeric_data, use = "complete.obs")
png("correlation_heatmap.png", width = 800, height = 600)
corrplot(correlation_matrix, method = "color", title = "Feature Correlations", tl.cex = 0.8)
dev.off()

# 2. Boxplot for Price Distribution
png("price_boxplot.png", width = 800, height = 600)
ggplot(data, aes(x = "", y = price_numeric)) +
  geom_boxplot(fill = "skyblue") +
  labs(title = "Price Distribution of Condominiums", y = "Price (MYR)", x = "") +
  theme_minimal()
dev.off()

# 3. Scatter Plot for Price vs. Property Size
png("price_vs_size.png", width = 800, height = 600)
ggplot(data, aes(x = property_size_numeric, y = price_numeric)) +
  geom_point(alpha = 0.6, color = "blue") +
  labs(title = "Price vs. Property Size", x = "Property Size (sq.ft.)", y = "Price (MYR)") +
  theme_minimal()
dev.off()
```

8. Commit: Initial commit
   Details: Created the repository and added basic structure to begin the project.
9. Commit: Added Readme Description
   Details: Added an initial draft of the README file, including an overview of the project and setup instructions.
10. Commit: Updated README.md
    Details: Enhanced the README with a more detailed project structure and contributions guidelines.
11. Commit: Added Houses.csv file
    Details: Uploaded the primary dataset (houses.csv) used for the analysis.
12. Commit: Added documentation for the analysis
    Details: Created detailed documentation for the project, including descriptions of visualizations and statistical methods.
13. Commit: R script for Data analysis
    Details: Uploaded R scripts for data preprocessing, visualization, and regression analysis.
14. Commit: Generated visualizations and analysis results
    Details: Created and saved visualizations, along with the analysis results, in the /outputs directory.