

IE4102 INDEPENDENT STUDY COURSE



Unsupervised and supervised learning on Singapore housing dataset

Submitted by

Low Eeron A0216716E

DEPARTMENT OF INDUSTRIAL SYSTEMS ENGINEERING AND

MANAGEMENT

NATIONAL UNIVERSITY OF SINGAPORE

Semester 1 2023

Summary

Introduction

The dynamic nature of the housing market necessitates the use of robust data analysis tools to aid potential buyers and sellers in making informed decisions. This report details the application of unsupervised and supervised machine learning techniques to the housing resale price dataset with a dual objective: developing a recommendation system and predicting housing resale prices.

Problem Definition and purpose

Housing decisions are often based on historical data, which sometimes may be incomplete or not representative of niche housing parameters. Buyers face the uncertainty of determining a fair price for a property, while sellers grapple with listing their homes at prices that reflect their true value. This project aims to address these challenges through:

- a. Unsupervised Learning: This approach is designed to help users input their desired housing parameters and receive recommendations based on houses that have sold in the past with similar attributes.
- b. Supervised Learning: This aspect of the project focuses on predicting a house's resale price based on varying parameters. It offers an estimation tool for users, providing insights into potential property values.

Key results

The culmination of this project resulted in the creation of two applications tailored to meet the needs of the housing market:

- a. A recommendation system built upon unsupervised learning techniques that offer insights into properties sold in the past based on user-specified parameters.
- b. A resale price prediction tool developed using supervised learning that provides users with a ballpark estimate, ensuring informed decision-making during transactions.

These applications ensure that both buyers and sellers have an upper hand in understanding and negotiating prices, bridging the existing information gap.

Conclusion

This project successfully marries unsupervised and supervised learning techniques to offer users a comprehensive perspective of the housing market. The recommendation system and the price prediction tool, when used in tandem, promise a more transparent and data-driven approach to housing transactions, benefiting all stakeholders involved.

Recommendation

For effective housing exploration, start with the recommendation system to find properties and make informed decisions. Then, use the resale price prediction tool for insights on transaction values. Keep datasets updated with quarterly resale prices and retrain models for precision. Stakeholders should prioritize regular tool updates. Improve prediction accuracy by adding diverse data, such as macroeconomic indicators and socio-demographic trends, for a comprehensive market view.

Acknowledgements

I express my sincere gratitude to Senior Lecturer Dr. Li Haobin for his invaluable guidance and unwavering support throughout the preparation of this report. Professor Li's insightful feedback and regular progress updates have been instrumental in shaping the quality of my work. I am truly grateful for the opportunity to benefit from his expertise.

Furthermore, I would like to extend my appreciation to the ISE department for granting me the opportunity to contribute to this project. The academic environment and resources provided have been crucial in the successful completion of this research endeavor.

A special acknowledgment goes to Teyang Lau, whose work served as a source of inspiration for my research in the supervised prediction housing resale part. Building upon Teyang Lau's approach, I chose to integrate the housing resale dataset with information on six key amenities to address the existing limitations in the depth of information available in the current housing resale price dataset. While our analyses share a common foundation, it's crucial to recognize the unique aspects of my work. The range of datasets used differs significantly; Teyang Lau utilized data spanning from 1990 to 2020, whereas I focused on Q2 of 2023 for my analysis. Additionally, I employed a combination of 1km and 2km radii in my predictive housing resale price analysis, and the parameters used for predicting housing resale prices were distinct, contributing to the distinctive nature of my research. [1]

Table of Contents

Contents

Summary	I
Acknowledgements	III
List of Figures	VII
1.0 Introduction	1
1.1 Problem definition.....	1
1.2 Scope	2
2.0 Literature Review	3
2.1 Distance Calculation Using distVincentySphere	3
2.2. Unsupervised Learning Techniques	4
2.2.1 K-means	4
2.2.2 Hierarchical Clustering	4
2.2.3 DBSCAN (Density-Based Spatial Clustering of Applications with Noise).....	5
2.3 Supervised Learning Techniques	5
2.3.1 Random Forest.....	5
2.3.2 XGBoost (Extreme Gradient Boosting)	6
2.3.3 Support Vector Machines (SVM)	6
2.3.4 Gradient Boosting Machine (GBM)	6
2.3.5 Lasso and Ridge Regression.....	7
2.4 Evaluation Metrics in Supervised Learning	7
2.4.1 Root Mean Square Error (RMSE)	7
2.4.2 Mean Absolute Percentage Error (MAPE).....	8
3.0 Methodology	9

3.1 Geospatial Data Handling and Proximity Analysis: An Overview	9
3.2 Data Acquisition and Refinement	10
3.2.1 Dataset Compilation and Selection.....	10
3.2.2 Geospatial Augmentation of the Base Dataset	10
3.2.3 Extraction of Spatial Coordinates for Amenities	10
3.2.4 Spatial Proximity Analysis	11
3.3 Data Cleaning and Enhancement	11
3.3.1 Data Integrity and Quality	11
3.3.2 Handling Missing Data	11
3.3.3 Data Validation	12
3.3.4 Data Refinement and Preliminary Analysis	12
3.4 Feature Engineering	12
3.4.1 Temporal Decomposition	12
3.4.2 Encoding and Spatial Feature Creation.....	12
3.5 Model Development, Validation, and Interpretation	13
3.5.1 Model Assumptions	13
3.5.2 Model Validation and Overfitting	13
3.5.3 Model Interpretation	13
3.6 Model Comparison.....	14
3.7 Comprehensive Data Management	14
3.8 Crafting the Shiny App	14
3.8.1 Recommendation System (Unsupervised Learning)	14
3.8.2 Predictive Models (Supervised Learning).....	15
3.9 Limitations and Challenges	15
4.0 Results.....	16

4.1 Preliminary Analysis of the Data	16
4.2 Model Selection Based on Performance Metrics	20
4.3 Feature importance on random forest and xgboost	22
4.4 Shiny app	24
4.4.1 Unsupervised Recommendation System.....	24
4.4.2 Supervised Predictive Housing resale price	26
5.0 Discussion	27
5.1 Synthesizing Model Outcomes	27
5.2 Spatial Distribution and Urban Dynamics	27
5.3 Integrating Analysis with User Experience	27
5.4 Implications and Recommendations	28
5.5 Limitations and Future Directions	28
6.0 Conclusion.....	29
7.0 Recommendations.....	30
7.1 Structured Housing Exploration	30
7.1.1 Initial Exploration with Recommendation System.....	30
7.1.2 Informed Decision Making through the Resale Price Prediction Tool	30
7.2 Continuous Data Updates and Training.....	30
7.3 Diversify Data Sources.....	31
List of References	32
Appendix A : Combining housing dataset with the 6 amenities	35
Appendix B: Preliminary analysis of the combined dataset	41
Appendix C: unsupervised and supervised analysis of the combined dataset	43

List of Figures

Figure 3.1: Overview of Data analysis	9
Figure 4.1: Average distances in km arrange by town	16
Figure 4.2: Average number of facilities arrange by town	18
Figure 4.3: Summary statistics of the amenities	19
Figure 4.4: Comparison of RMSE across Different Models	21
Figure 4.5: Comparison of MAPE across Different Models	21
Figure 4.6: Random forest feature importance	22
Figure 4.7: Xgboost feature importance	23
Figure 4.8: Housing recommendation part 1	24
Figure 4.9: Housing recommendation part 2	25
Figure 4.10: Predicting housing resale price	26
Figure A.1: Adding longitude and latitude to original Q2 2023	35
Figure A.2: Getting parks and hawkers longitude and latitudes	36
Figure A.3: Adding number of hawkers within 1km radius , and the closest hawker distance	36
Figure A.4: Adding number of mrt within 1km radius , and the closest mrt distance	37
Figure A.5: Adding number of park within 1km radius , and the closest park distance	37
Figure A.6: Adding longitude and latitude in the school dataset from the postal code	38

Figure A.7: Adding number of school within 2km radius , and the closest school distance	38
Figure A.8: Adding number of shopping mall within 1km radius , and the closest shopping mall distance	39
Figure A.9: Adding longitude and latitude in the supermarket dataset from the postal code	39
Figure A.10: Adding number of supermarket within 1km radius, and closest supermarket distance	40
Figure B.1: Data cleaning to include the relevant columns for future data analysis	41
Figure B.2: Preliminary analysis part 1	42
Figure B.3: Preliminary analysis part 2	42
Figure C.1: unsupervised recommendation system shiny app part 1	43
Figure C.2: unsupervised recommendation system shiny app part 2	44
Figure C.3: unsupervised recommendation system shiny app part 3	45
Figure C.4: unsupervised recommendation system shiny app part 4	46
Figure C.5: unsupervised recommendation system shiny app part 5	47
Figure C.6: unsupervised recommendation system shiny app part 6	47
Figure C.7: Training and testing model part 1	48
Figure C.8: Training and testing model part 2	49

Figure C.9: Training and testing model part 3	49
Figure C.10: Manually predicting housing resale price based on the parameters	50
Figure C.11: Supervised predicting housing resale price shiny app	50

1.0 Introduction

1.1 Problem definition

The modern housing market in Singapore is characterized by its inherent fluidity, inundating potential homeowners and renters with vast amounts of information. Platforms like PropertyGuru offer basic filtering mechanisms, such as property type, price, number of bedrooms, and even by MRT stations, districts, and HDB estates[2]. Yet, these often fall short of discerning buyers' comprehensive needs. While some prioritize proximity to essential amenities, others may seek a more concentrated presence of specific amenities. This underscores a gap in meeting the diverse needs of a varied user base.

Compounding the issue are observed price discrepancies on platforms like PropertyGuru, where similar properties in the same area may have varying listings[2]. This inconsistency complicates decision-making and amplifies the need for better data transparency.

Beyond property discovery lies the challenge of interpreting abundant data. With a plethora of listings and information, stakeholders face the risk of overpricing or undervaluing properties. Determining a property's true market value in a dynamic environment like Singapore becomes an intricate endeavor.

Given this backdrop, two primary challenges emerge:

1. The complex task of navigating vast listings, with limited filtering options that may not fully address the unique needs of many buyers.
2. The essential need for accurate interpretation of data to ensure optimal property pricing in a fast-evolving market.

1.2 Scope

Central to this research is the application of machine learning, a powerful tool that can uncover patterns and insights from vast datasets, effectively addressing the aforementioned challenges. The goals are twofold:

Data Structuring for Tailored Insights: Aimed at simplifying the property search process, this facet of the research focuses on structuring and categorizing housing data, ensuring stakeholders can efficiently identify properties meeting their criteria.

Predictive Insights on Housing Resale Values: To address transaction pricing challenges, the research emphasizes accurate market value forecasting. In a market influenced by countless variables, predictive tools become vital.

By converging these objectives, this research offers a comprehensive solution: streamlining property searches and offering foresight into potential market values.

2.0 Literature Review

As the intricacies of housing market analysis are delved into, it is reminded that every discovery and insight stands on the shoulders of preceding work. A spotlight is cast on the foundational knowledge and tools that have shaped the field of geospatial and housing market analytics by the subsequent literature review. These established techniques, borne out of rigorous academic inquiry and practical application, have paved the way for refined approaches to data analysis. Once a comprehensive understanding of the literature is attained, an in-depth exploration will follow, focusing on the specific methodologies tailored for our study. Each methodology has been chosen for its precision, relevance, and efficacy in deciphering complex datasets like ours.

2.1 Distance Calculation Using distVincentySphere

In the realm of geospatial analysis, accurately computing the distance between two points on the Earth's surface becomes paramount. The Earth's shape, deviating from a perfect sphere to resemble an oblate spheroid, complicates such calculations. Nevertheless, for many applications, where ultra-high precision isn't the end goal or over distances where the Earth's subtle flattening is less impactful, approximating Earth as a sphere proves reasonable. This leads to the use of the distVincentySphere, a derivative of Vincenty's formula that assumes a spherical Earth. This method finds favor for its commendable balance between precision and computational simplicity, especially in scenarios prioritizing performance [3].

In the context of the Singapore housing market, geospatial analysis is considered indispensable. The value of a property can be significantly influenced by its proximity to amenities, transportation hubs, and other urban landmarks. The distVincentySphere is leveraged so that distances between properties and these critical points of interest can be

accurately calculated, ensuring that spatial attributes of properties are correctly captured and represented. Furthermore, it is recognized that, given Singapore's compact urban landscape, notable value differences can be led to by even minor distance variations, highlighting the essential need for accurate distance measurements.

2.2. Unsupervised Learning Techniques

Unsupervised learning excels in situations where the goal is to identify inherent patterns or structures in data without explicit labels [4]. Given the multifaceted attributes of the housing dataset, these techniques can be instrumental in uncovering hidden relationships or clusters.

Clustering aids in categorizing data into groups based on similarity measures. Such categorization can be pivotal for datasets like the one at hand, which encompasses various attributes ranging from flat type to proximity to amenities [5].

2.2.1 K-means

Segmenting the housing dataset into 'K' clusters based on similarity, K-means can discern common housing patterns or categories. The algorithm minimizes the intra-cluster variance, ensuring that data points in the same cluster are as similar as possible [6].

Within the housing dataset, K-means can be instrumental in segmenting properties into distinct clusters based on attributes like size, location, and type. This can offer potential buyers a more tailored list of properties that fit within a specific category of interest

2.2.2 Hierarchical Clustering

This method provides a tree-like structure of clusters, making it suitable for understanding data at various granularities. Hierarchical clustering is especially pertinent when there's potential nested structure in the data, such as housing categories and subcategories [7].

Given the multi-layered nature of housing attributes, hierarchical clustering can aid in understanding how properties can be categorized at multiple levels. For instance, the first level might cluster based on property type, while subsequent levels could focus on price brackets or proximity to amenities

2.2.3 DBSCAN (Density-Based Spatial Clustering of Applications with Noise)

Unlike K-means, which assumes clusters to be spherical, DBSCAN can detect clusters of arbitrary shapes. This is particularly useful for complex datasets where clusters may not be clearly defined. Moreover, its ability to distinguish noise or outliers ensures that anomalous data points don't influence the clustering unduly [8].

For a dataset as varied as housing, where properties might not always fit neatly into spherical clusters, DBSCAN provides flexibility. It can highlight unique or niche property groupings that other algorithms might miss, ensuring a more comprehensive analysis

2.3 Supervised Learning Techniques

In scenarios with labeled data, such as housing datasets with known resale prices, supervised learning techniques become indispensable. These techniques model the relationship between input attributes and the target variable. In the context of our study, supervised learning is employed to predict housing resale prices based on various features of the data [9].

2.3.1 Random Forest

An ensemble method, Random Forest constructs multiple decision trees, enhancing prediction accuracy and robustness against overfitting. This technique can handle a wide range of data attributes without extensive preprocessing [10].

In predicting housing prices, the Random Forest algorithm can accommodate the myriad of features associated with properties, from basic attributes like size and age to more nuanced ones like view quality or architectural significance

2.3.2 XGBoost (Extreme Gradient Boosting)

Celebrated for its computational efficiency and prediction accuracy, XGBoost is a scalable take on gradient boosting. It builds trees in a sequential manner, where each tree corrects the errors of its predecessor [11].

With its ability to iteratively enhance prediction accuracy, XGBoost can refine housing price predictions, ensuring they resonate with the dynamic Singapore housing market trends

2.3.3 Support Vector Machines (SVM)

SVMs are adept at both classification and regression tasks. They work by identifying an optimal hyperplane that best separates data into classes. The kernel trick in SVM allows it to handle non-linear data, providing flexibility in modeling complex relationships [12].

SVM's flexibility in handling non-linear relationships makes it a strong contender for modeling the intricate variables influencing housing prices

2.3.4 Gradient Boosting Machine (GBM)

GBM, like XGBoost, builds trees iteratively. Each tree aims to correct the residuals or errors from the preceding trees, ensuring that prediction accuracy improves with each iteration [13].

Much like XGBoost, GBM's iterative approach can be instrumental in fine-tuning predictions, ensuring that as the model learns, it gets progressively closer to actual market values

2.3.5 Lasso and Ridge Regression

Both are regularization techniques for linear regression. Lasso performs feature selection, ensuring irrelevant features don't influence the model, while Ridge provides robustness against multicollinearity, ensuring stable coefficient estimates [14,15].

With the multitude of potential influencing factors in a housing dataset, Lasso can ensure that only the most relevant features impact the model. Simultaneously, Ridge Regression provides a safeguard against any inter-correlation between features, ensuring a more stable and reliable prediction

2.4 Evaluation Metrics in Supervised Learning

Evaluation metrics ensure that the supervised models' predictions align well with the actual outcomes, providing a quantitative measure of performance.

2.4.1 Root Mean Square Error (RMSE)

RMSE represents the square root of the average squared differences between the predicted and actual values. In simpler terms, it measures the average magnitude of errors made by a predictive model [16].

In the context of predicting housing resale prices, RMSE provides a clear indication of the model's overall prediction accuracy in terms of the actual monetary difference. A lower RMSE indicates that the model's predictions are closer to the actual resale prices, making it a vital metric for gauging the model's effectiveness. Given the substantial financial implications of housing transactions, even small prediction errors can translate to significant monetary discrepancies. Hence, minimizing RMSE becomes paramount to ensure stakeholders make informed decisions based on the model's outputs.

2.4.2 Mean Absolute Percentage Error (MAPE)

MAPE calculates the average of the absolute percentage differences between the predicted and actual values. It provides a relative measure of error, expressing prediction inaccuracies as a percentage [17].

While RMSE provides an absolute measure of error in terms of monetary units, MAPE offers a relative perspective, making it easier for stakeholders to contextualize the model's prediction accuracy. For instance, a MAPE of 5% suggests that, on average, the model's predictions deviate from actual prices by 5%. This percentage representation can be especially intuitive for stakeholders, allowing them to gauge the model's reliability in a more comprehensible format. In the dynamic housing market of Singapore, where prices can vary widely across districts and property types, having a relative measure like MAPE can be particularly insightful for potential buyers, sellers, and investors.

3.0 Methodology

An in-depth exploration into Singapore's housing resale market was undertaken, blending data science, spatial analytics, and machine learning techniques. This section elucidates the intricate methodologies employed.

3.1 Geospatial Data Handling and Proximity Analysis: An Overview

To provide a brief overview of our research approach: Datasets encompassing housing resale data and various amenities were primarily acquired from Data.Gov.sg and Kaggle. Missing spatial data were supplemented using geocoding, converting addresses or postal codes to latitude and longitude details. Subsequent spatial analyses determined property proximity to amenities and the density of amenities within specific radii. Following this geospatial integration, a recommendation system was developed, and predictive models for housing resale prices were formulated. Figure 3.1 offers a visual summary of this workflow.

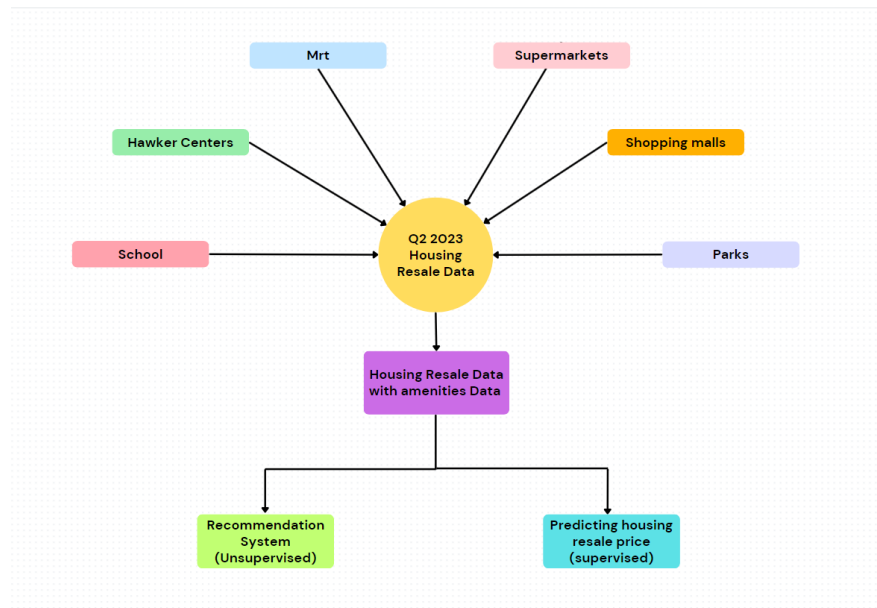


Figure 3.1: Overview of Data analysis

With this visual and descriptive guide, the subsequent sections will delve into the detailed methodologies and their nuances.

3.2 Data Acquisition and Refinement

3.2.1 Dataset Compilation and Selection

Data for the study was sourced from Data.Gov.sg [18], a recognized data repository. The Q2 2023 dataset was chosen to offer a contemporary perspective on housing resale trends. The research incorporated datasets including: Housing Resale Data [18], School Data [19], Hawker Centers [20], MRT Stations [21], Supermarkets [22], Shopping Malls [23], Parks [24]

3.2.2 Geospatial Augmentation of the Base Dataset

The base dataset was initially devoid of latitude and longitude details. To address this, the ggmap package was utilized in tandem with the Google Maps API to geocode addresses, thereby equipping each property with accurate spatial coordinates.

3.2.3 Extraction of Spatial Coordinates for Amenities

For parks and hawker centers, KML files were processed to retrieve their spatial coordinates. Using the `st_read` function, KML data was read, and coordinates were subsequently extracted. For schools and supermarkets, where latitude and longitude were absent, a dedicated function leveraging the ggmap package and Google Maps API was crafted to geocode postal codes. This function transformed postal codes into accurate spatial markers.

3.2.4 Spatial Proximity Analysis

Once all datasets were enriched with spatial coordinates (longitude and latitude), a comprehensive spatial analysis was conducted. The proximity of each property to various amenities was determined:

For all amenities, including hawker centers, MRT stations, parks, shopping malls, supermarkets, and even schools, two primary metrics were derived: distance to the closest facility and number of facilities within a specific radius. However, while a 1km vicinity was considered for most amenities, schools were evaluated within a broader 2km boundary due to their distinct importance in housing decisions. The decision to consider a 2km radius for schools was informed by the understanding that a larger catchment area increases the likelihood of children gaining admission. This expanded radius is deemed beneficial for families with school-going children, as proximity to preferred educational institutions often plays a pivotal role in housing decisions [25].

3.3 Data Cleaning and Enhancement

3.3.1 Data Integrity and Quality

For coherent analyses, columns with mixed temporal data, like 'remaining_lease', were split into separate 'years' and 'months' columns. Further, categorical variables, including 'town', 'flat_type', 'storey_range', and 'flat_model', were transformed via one-hot encoding, making them amenable to machine learning techniques.

3.3.2 Handling Missing Data

In cases where certain amenities lacked spatial coordinates, innovative geocoding methods were employed. Postal codes available for these amenities were transformed into latitude and longitude markers using the aforementioned geocoding function.

3.3.3 Data Validation

Post data acquisition, a validation phase was incorporated. The constructed price prediction model was tested against the latest Q3 2023 data, and the predictions were found to be within Q1 and Q3 of the predicted range, affirming the model's reliability.

3.3.4 Data Refinement and Preliminary Analysis

Upon validation, our study embarked on thorough data cleaning, eliminating non-compliant entries and superfluous columns, enhancing data reliability. This refined data was cataloged in a 'final' CSV file, earmarked for ensuing inquiries. In tandem, insightful data visualizations, encompassing bar charts were crafted, offering an instant, lucid comprehension of the geographical spread of amenities. This phase also embraced a focused descriptive statistical analysis, involving key summary statistics — mean, median, and others — to capture the data's essence. These processes not only affirmed the data's robustness but also set the stage for in-depth, successive analyses.

3.4 Feature Engineering

3.4.1 Temporal Decomposition

The 'remaining_lease' column, teeming with combined temporal data of years and sometimes months, was elegantly decomposed to facilitate a more granulated temporal scrutiny.

3.4.2 Encoding and Spatial Feature Creation

Ensuring the dataset's compatibility with machine learning models, categorical variables underwent a transformation via one-hot encoding. Moreover, the spatial coordinates obtained through geocoding were instrumental in devising spatial features. These features

encompassed proximity metrics and density measurements, which added depth to the dataset's spatial narrative.

3.5 Model Development, Validation, and Interpretation

3.5.1 Model Assumptions

Several assumptions were made during the modeling phase. It was assumed that external monetary conditions had no notable effect on the resale price fluctuations. Furthermore, popular amenities like renowned schools, shopping malls, or MRT stations were assumed not to exert a significant impact on the resale prices. Given that the dataset combines data from different years, it was assumed that such combinations did not introduce crucial disparities in the results or the model's performance.

3.5.2 Model Validation and Overfitting

To fathom the performance of models beyond their training data, the dataset was bifurcated into training and test subsets. Embracing cross-validation techniques fortified the models against overfitting, ensuring their robustness. Regularization techniques, including LASSO and Ridge regression, were integrated, placing constraints on coefficient magnitudes and ensuring model generalizability. For hyperparameter tuning, a conscious decision was made not to overfit the model. The Mean Absolute Percentage Error (MAPE) was observed to be less than 5% for both Random Forest and XGBoost, indicating satisfactory model performance.

3.5.3 Model Interpretation

The Random Forest model, an ensemble methodology, was trained, unveiling feature importance and providing a spectrum of expected resale prices. Concurrently, the XGBoost model was leveraged to capture intricate non-linear relationships, highlighting influential

variables. The robustness of models was further validated by computing the Root Mean Square Error (RMSE) for an array of models, ensuring a comprehensive and rounded analysis.

3.6 Model Comparison

An array of models was trained, and their performance was rigorously compared using RMSE and MAPE. The Random Forest and XGBoost models were ultimately chosen due to their superior performance metrics.

3.7 Comprehensive Data Management

For the sanctity of consistency and future replication potential, an unwavering data management protocol was adopted. Every significant analytical stride and its results were meticulously documented and archived in a structured CSV format, laying the groundwork for subsequent analyses or visualizations.

3.8 Crafting the Shiny App

At the heart of the digital solution is a Shiny app, a product of meticulous design and tailored functionalities, aimed specifically at potential homeowners in Singapore.

3.8.1 Recommendation System (Unsupervised Learning)

Dynamic User Interface: The Shiny app's recommendation system offers a dynamic and interactive interface where users can specify their preferences, from proximity to amenities to types of neighborhoods. This interface is not just about input; it's about experience. It offers users a chance to navigate through various clustering methods, allowing for customization at every step. The culmination of this is a detailed map, visually representing recommended housing areas based on user-defined criteria.

Behind the user interface, a sophisticated server logic operates, processing user inputs, filtering data accordingly, and applying the appropriate clustering techniques. Once clusters have been established, specific clusters are highlighted based on predefined parameters and clustering criteria, offering users a curated list of recommendations.

3.8.2 Predictive Models (Supervised Learning)

Data Preparation: Before any predictive modeling could commence, the dataset underwent thorough preprocessing. Essential transformations were applied to ensure that the data was in the optimal format for modeling. It was then divided into training and test subsets, setting the stage for the subsequent model training phase.

Model Training and Evaluation: An array of regression models was trained on the dataset. Each model's performance was rigorously evaluated using established metrics. This comparative analysis allows for the identification of the most accurate models for price prediction.

In addition, a Manual Prediction Shiny App was introduced. Here, users can input specific housing attributes. These inputs are then processed and used by the trained models to predict resale prices. A buffer is added to account for uncertainties, ensuring users receive a prediction range, offering a more realistic and practical forecast.

3.9 Limitations and Challenges

One notable limitation of this study was the dataset's recency. Although efforts were made to incorporate the most recent data, the dataset for amenities wasn't the absolute latest available. Another challenge was the amalgamation of data across different years, which brought its own set of complexities.

4.0 Results

4.1 Preliminary Analysis of the Data

An initial exploration of the dataset was conducted, providing insights into the spatial distribution of amenities in relation to housing units across different towns in Singapore. The findings, illustrated by the provided visualizations and summary table, offer a nuanced perspective on the housing landscape in Singapore.

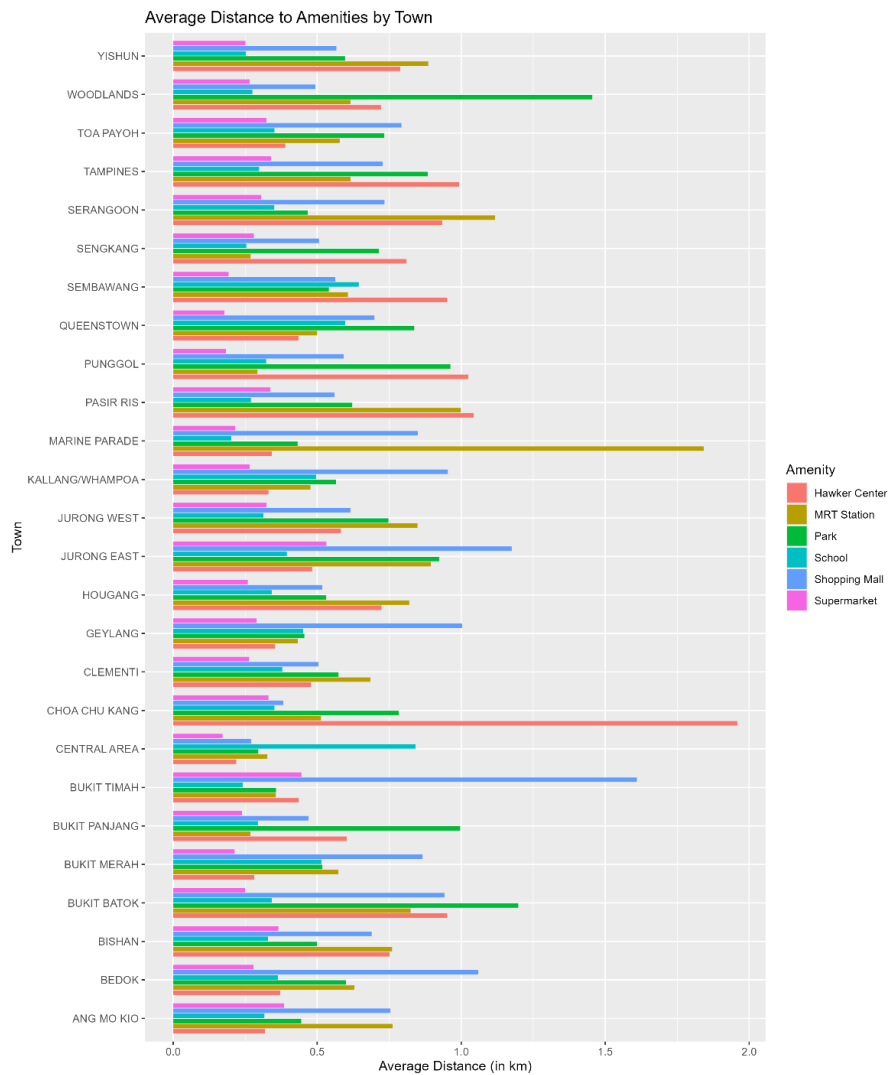


Figure 4.1: Average distances in km arrange by town

Figure 4.1 delineates the average distances from housing units to various amenities across different towns. Predominantly, it can be observed that the majority of the towns have amenities situated within an average distance of 1km. However, there are specific exceptions that stand out: Woodlands residents, on average, find themselves approximately 1.5km away from the nearest park. In Marine Parade, the average distance to the closest MRT station is about 1.75km. Choa Chu Kang exhibits a slightly extended average distance of around 2km to the nearest hawker center. For those residing in Bukit Timah, the average distance to a shopping mall is about 1.6km.

These variations emphasize the spatial heterogeneity in amenity accessibility across different towns and underscore the importance of geospatial analysis in understanding urban dynamics.

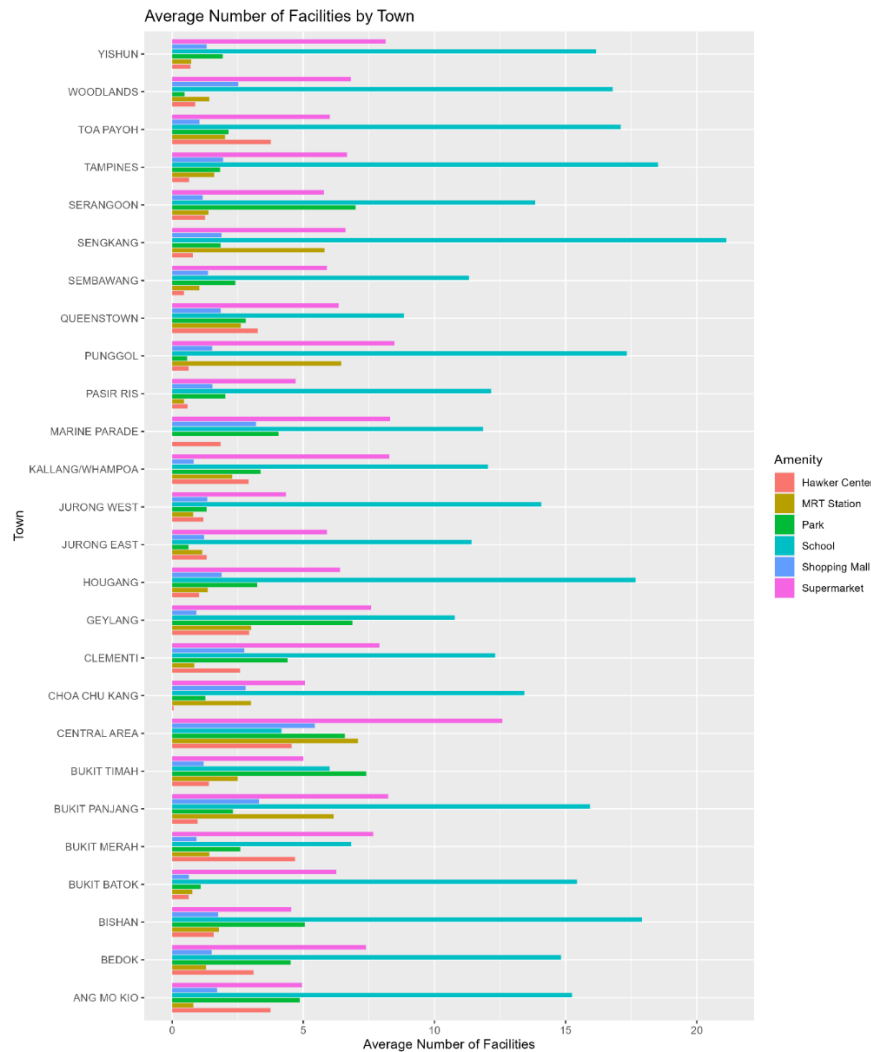


Figure 4.2: Average number of facilities arrange by town

Figure 4.2 illustrates the density of various amenities across different towns. Overall, towns seem to exhibit a consistent number of facilities, indicating a balanced distribution of resources. Specific observations that stand out are:

Central Area displays a relatively lower density of school facilities compared to other regions, suggesting potential urban planning strategies that prioritize other amenities. Sengkang stands out with the highest density of schools, reflecting the area's potential appeal to families with school-going children. When it comes to supermarkets, the Central

Area leads with the highest density, catering to the urban population's daily needs. On the contrary, Jurong West appears to have a slightly lower density of supermarkets.

These observations underscore the unique amenity profiles of each town, with each catering to distinct resident needs and urban planning priorities.

Descriptive Statistics of Distances and Counts:

	Min	Median	Mean	Max	X1st.Qu.,25.	X3rd.Qu.,75.
<i>hawker_closest_distance</i>	0.033691061138632	0.840130826777672	0.761725473886178	2.80398433388695	0.371057398578678	1.02172033451958
<i>num_hawker_1km</i>	0	1	1.43161343161343	9	0	2
<i>mrt_closest_dist</i>	0.0264104881460214	0.552227521337943	0.617313454533644	3.41317620024466	0.336933801516211	0.825955311936475
<i>num_mrt_1km</i>	0	1	2.38886158886159	14	1	3
<i>park_closest_dist</i>	0.0681389534230855	0.686109097951199	0.78765783301917	2.40322990458056	0.46813407138607	0.995607483107017
<i>num_park_1km</i>	0	1	2.28058968058968	15	1	3
<i>sch_closest_dist</i>	0.00812176886117962	0.293583337370964	0.345709935582443	3.29578401975786	0.195406361779354	0.426793142889845
<i>num_sch_2km</i>	0	16	15.3772317772318	29	12	18
<i>shoppingmall_closest_dist</i>	0.00505698397679216	0.595793485996857	0.66374383359772	3.11757249627717	0.391263429374416	0.873920675335062
<i>num_shoppingmall_1km</i>	0	2	1.6972972972973	15	1	2
<i>supermarket_closest_dist</i>	0	0.251493000806214	0.276753460509932	3.09937827460337	0.161356218160154	0.367957784724172
<i>num_supermarket_1km</i>	0	6	6.69041769041769	18	5	9

Figure 4.3: Summary statistics of the amenities

Figure 4.3 above provides a detailed perspective on the dataset's metrics, offering insights into the central tendencies and spread of distances to amenities and their respective counts. Specific observations from the table include:

Supermarkets stand out with the lowest average distance from housing units, signifying their widespread presence and accessibility in Singapore. This is further emphasized by the fact that all amenities, on average, are located within a 1km radius of housing units, underscoring the well-planned urban infrastructure of Singapore. Contrarily, parks exhibit the highest average distance, suggesting that while they are integral to Singaporean neighborhoods, they might be spaced out a bit more compared to other amenities.

When considering the density of amenities within a 1km radius, supermarkets again lead with the highest average count, emphasizing their prevalence in residential areas. On the other hand, hawker centers, despite being iconic to Singapore's culture, have the lowest average count within the same radius.

These metrics elucidate the intricate balance of amenities in Singapore, highlighting the emphasis on both accessibility and distribution in urban planning. In light of these findings, it was concluded that each town in Singapore offers a unique blend of accessibility and amenity density. These preliminary insights pave the way for more in-depth analyses in the subsequent sections of the report.

4.2 Model Selection Based on Performance Metrics

Based on the performance metrics, particularly the Root Mean Square Error (RMSE) and Mean Absolute Percentage Error (MAPE), Random Forest and XGBoost were chosen as the primary models for predicting resale prices of housing units.

The selection of these models is evident when observing the RMSE and MAPE figures. Both metrics provide insight into the accuracy of predictions made by the models. Lower RMSE values indicate better fit to the data, while lower MAPE values suggest the model's predictions are, on average, closer in percentage terms to the actual values.

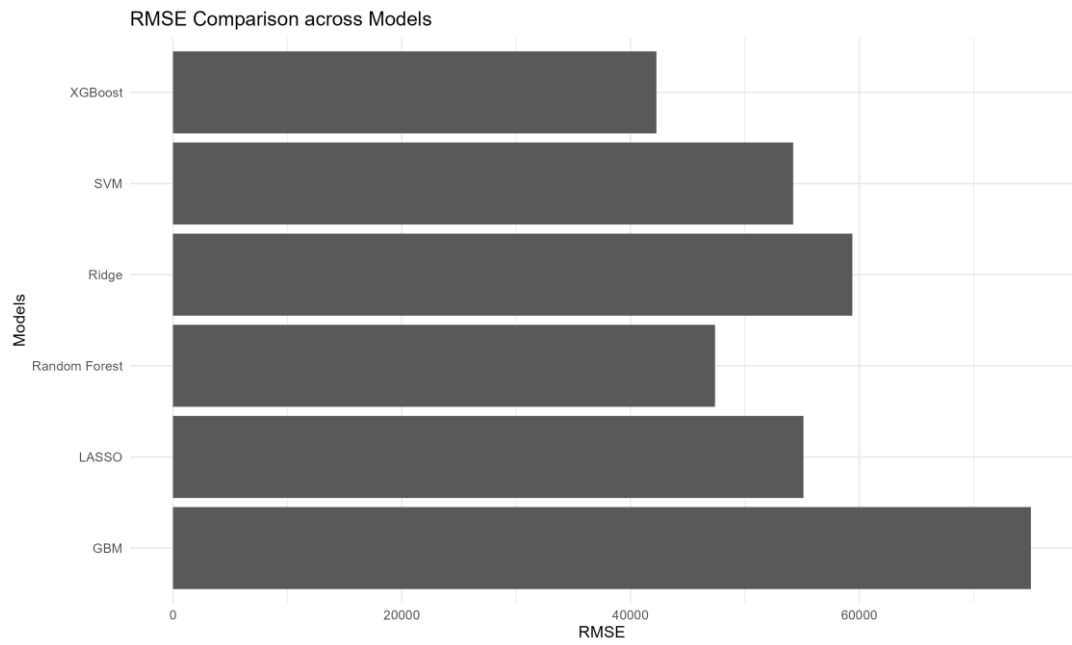


Figure 4.4: Comparison of RMSE across Different Models

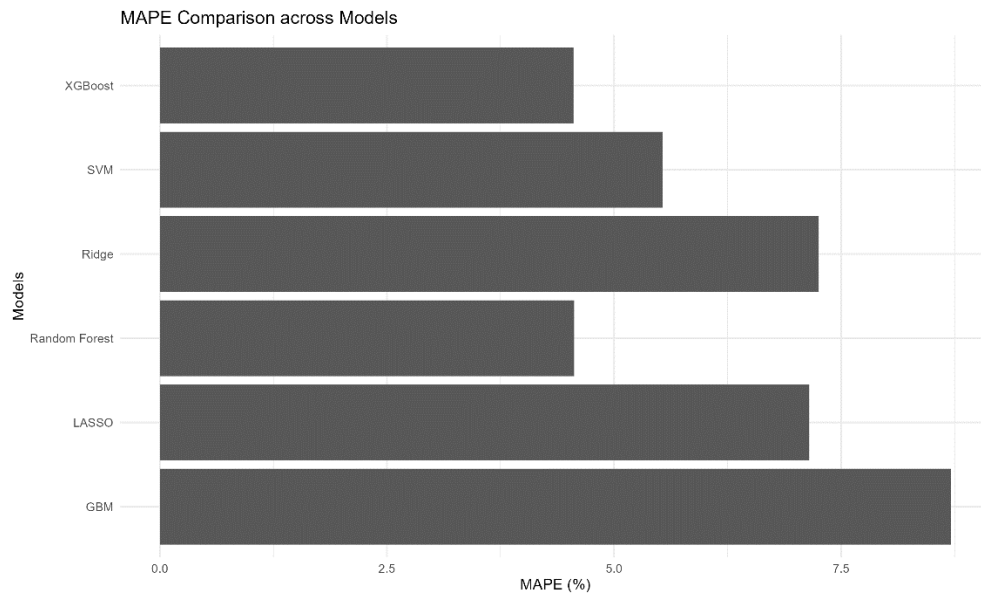


Figure 4.5: Comparison of MAPE across Different Models

From Figures 4.4 and 4.5, it's evident that both Random Forest and XGBoost outperformed the other models in terms of these metrics, making them the ideal choice for this dataset and the problem at hand.

4.3 Feature importance on random forest and xgboost

The significance of features, as determined by machine learning models, can provide profound insights into the factors that most influence the resale prices of housing units.

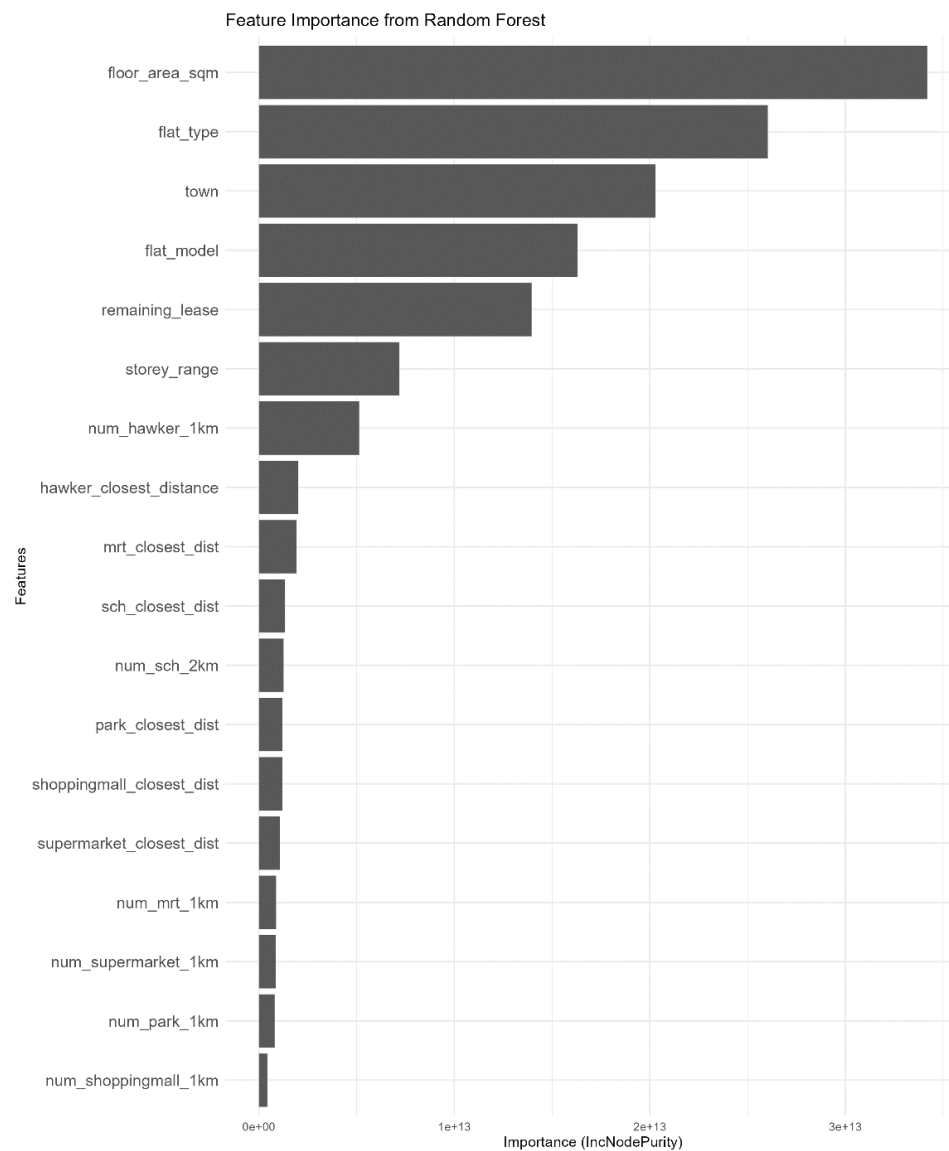


Figure 4.6: Random forest feature importance

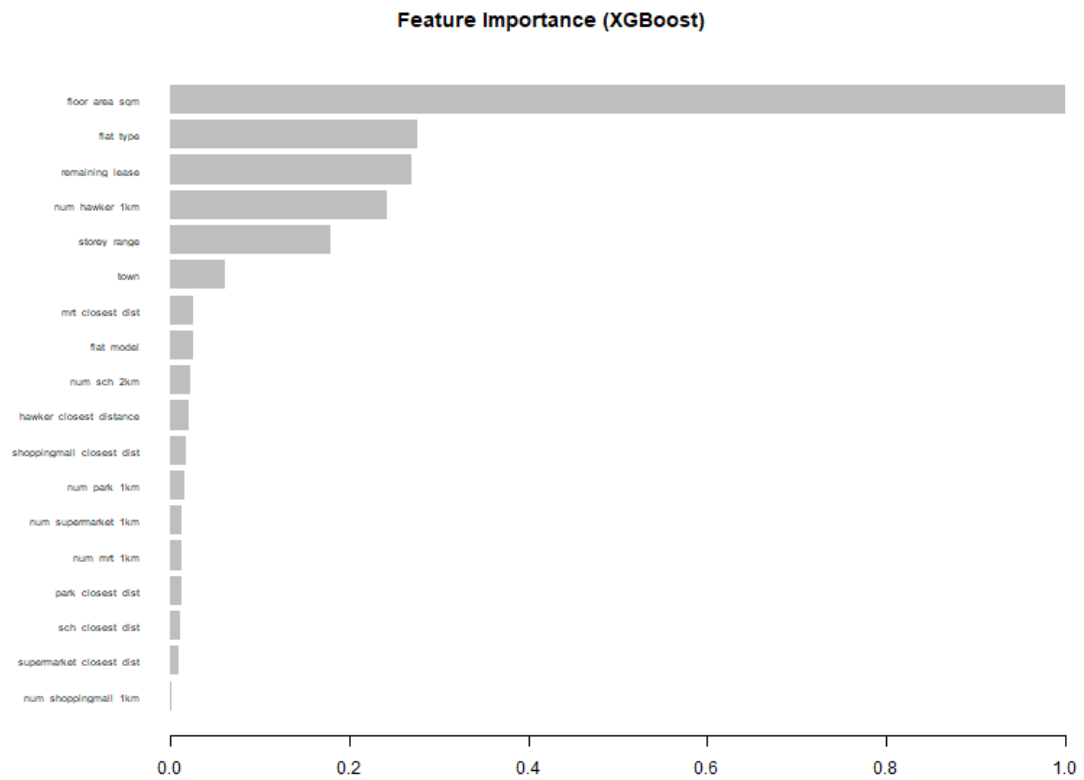


Figure 4.7: XGBoost feature importance

Figure 4.6 Random Forest reveals the model's priority on 'floor area', 'flat type', and 'town' in influencing resale prices, highlighting the critical roles of unit size, type, and location. Conversely, Figure 4.7 shows XGBoost's emphasis on 'floor area', 'flat type', and the unique 'remaining lease', pointing to the lease duration's unexpected significance. Despite these differences, both models concur on the paramountcy of 'floor area' and 'flat type'. This consistency across two distinct models solidifies the status of these features as reliable, universal predictors of resale value, regardless of the analytical method.

4.4 Shiny app

4.4.1 Unsupervised Recommendation System

The Shiny application, designed as a comprehensive tool for housing in Singapore, encompasses both an unsupervised recommendation system and a supervised predictive model for housing resale prices. The unsupervised recommendation system offers users a dynamic interface, allowing them to tailor their housing searches based on specific criteria, and subsequently receive a curated list of recommendations.

Housing Recommendations

Select Parameters:

- ☒ MRT
- ☐ Park
- ☐ School
- ☐ Supermarket
- ☐ Shopping Mall
- ☐ Hawker Center
- ☐ Town
- ☐ Flat Type
- ☐ Storey Range
- ☐ Flat Model
- ☐ Floor Area
- ☐ Remaining Lease

Select Clustering Method:

- ☒ K-means Clustering
- ☐ Hierarchical Clustering
- ☐ DBSCAN

Number of Clusters:

2

At most MRT Distance (km):

1

Minimum Number of MRT Stations (1Km) radius:

1

At most Park Distance (km):

1

Minimum Number of Parks (1Km) radius:

1

At most School Distance (km):

1

Minimum Number of Schools (1Km) radius:

month	town	flat_type	block	street_name	storey_range	floor_area_sqm	flat_model	lease_commence_date	remaining_lease	resale_price	Latitude	Longitude	hawker_closest_distance	num_h
2023-04	BUKIT BATOK	2 ROOM	450A	BT BATOK WEST AVE 6	10 TO 12	38.00	Model A	2019	94 years 11 months	320000.00	1.35	103.74		0.33
2023-04	BUKIT BATOK	2 ROOM	450A	BT BATOK WEST AVE 6	19 TO 21	38.00	Model A	2019	94 years 11 months	320000.00	1.35	103.74		0.33
2023-04	BUKIT BATOK	2 ROOM	292A	BT BATOK EAST AVE 6	04 TO 06	47.00	Model A	2019	94 years 10 months	340000.00	1.34	103.75		1.84
2023-04	BUKIT BATOK	2 ROOM	450A	BT BATOK WEST AVE 6	13 TO 15	47.00	2-room	2019	94 years 11 months	350000.00	1.35	103.74		0.33
2023-04	BUKIT BATOK	2 ROOM	450A	BT BATOK WEST AVE 6	19 TO 21	47.00	Model A	2019	94 years 11 months	350000.00	1.35	103.74		0.33
2023-04	BUKIT BATOK	3 ROOM	208	BT BATOK ST 21	04 TO 06	73.00	Model A	1983	59 years 05 months	388000.00	1.35	103.75		1.26
2023-04	BUKIT BATOK	3 ROOM	293D	BT BATOK ST 21	04 TO 06	68.00	Model A	2019	94 years 09 months	470000.00	1.35	103.76		1.88
2023-04	BUKIT BATOK	3 ROOM	450D	BT BATOK WEST AVE 6	07 TO 09	67.00	Model A	2019	94 years 11 months	475000.00	1.35	103.74		0.36
2023-04	BUKIT BATOK	3 ROOM	450A	BT BATOK WEST AVE 6	10 TO 12	67.00	Model A	2019	95 years	485000.00	1.35	103.74		0.33
2023-04	BUKIT BATOK	3 ROOM	296B	BT BATOK ST 22	10 TO 12	68.00	Premium Apartment	2018	93 years 11 months	490000.00	1.34	103.75		1.09
2023-04	BUKIT BATOK	4 ROOM	201	BT BATOK ST 21	01 TO 03	109.00	Model A	1985	61 years 03 months	498000.00	1.35	103.75		1.16
2023-04	BUKIT BATOK	3 ROOM	296B	BT BATOK ST 22	28 TO 30	68.00	Premium Apartment	2018	93 years 11 months	500000.00	1.34	103.75		1.09
2023-04	BUKIT BATOK	4 ROOM	211	BT BATOK ST 21	10 TO 12	104.00	Model A	1985	61 years 01 month	505000.00	1.35	103.75		1.37
2023-04	BUKIT BATOK	3 ROOM	450D	BT BATOK WEST AVE 6	19 TO 21	67.00	Model A	2019	95 years	505000.00	1.35	103.74		0.36
2023-04	BUKIT BATOK	4 ROOM	211	BT BATOK ST 21	10 TO 12	104.00	Model A	1985	61 years	510000.00	1.35	103.75		1.37
2023-04	BUKIT BATOK	4 ROOM	207	BT BATOK ST 21	07 TO 09	104.00	Model A	1983	59 years 05 months	513888.00	1.35	103.75		1.30
2023-04	BUKIT BATOK	4 ROOM	108	BT BATOK WEST AVE 6	04 TO 06	104.00	Model A	1985	61 years 01 month	520000.00	1.35	103.75		1.01
2023-04	BUKIT BATOK	4 ROOM	217	BT BATOK ST 21	10 TO 12	93.00	New Generation	1984	60 years 08 months	528500.00	1.35	103.75		1.70

Figure 4.8: Housing recommendation part 1

The development of an interactive recommendation system, as depicted in Figure 4.8, provides users with a sophisticated and user-friendly interface to customize their housing search based on specific criteria.

The system offers flexibility, allowing users to modify various parameters. They can fine-tune the clustering methods, determine the number of clusters, and even select the basis for clustering criteria.

Post customization, the interface presents a curated list of the top 45 housing recommendations on the right side of the shiny app. This list is tailored based on the user's selected parameters, ensuring that the recommendations align closely with their preferences.

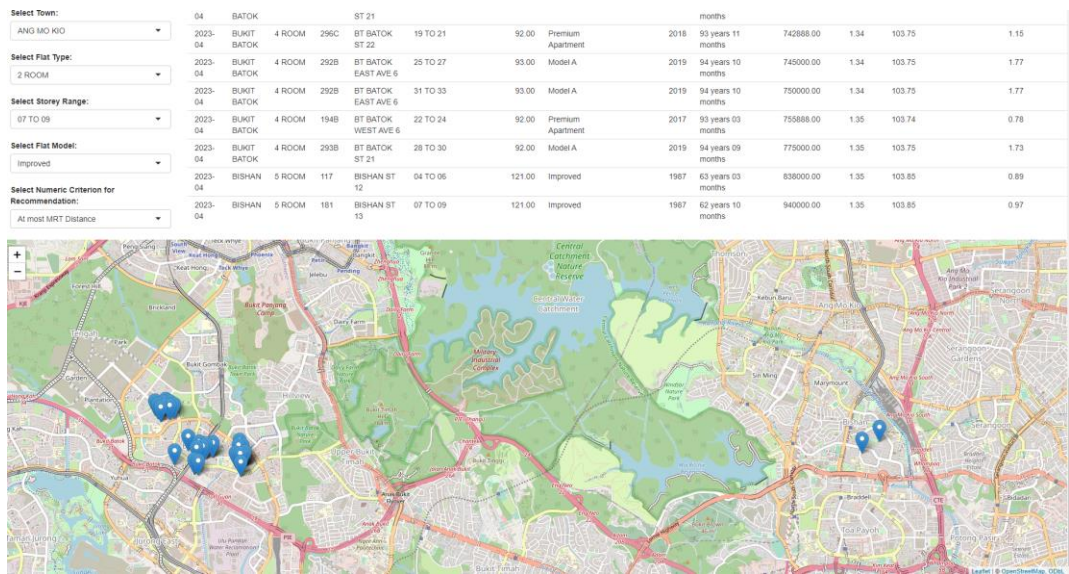


Figure 4.9: Housing recommendation part 2

Figure 4.9 displays an integrated map of Singapore at the bottom of the shiny app to enhance user experience and spatial understanding. This visual aid allows users to geographically situate the recommended housing units, aiding in decision-making by providing a spatial context.

This interactive recommendation system not only streamlines the housing search process but also empowers users by offering them control over the search criteria, ensuring a more personalized and efficient search experience.

4.4.2 Supervised Predictive Housing resale price

On the other hand, the supervised predictive model provides insights into potential housing resale prices based on various influential factors. Together, these components present users with a holistic view, assisting them in making informed housing decisions in the Singaporean market.

Predict Resale Price

Town: AND KIO KIO

Flat Type: 1 ROOM

Flat Model: 3Gen

Storey Range: 01 TO 03

Floor Area (sqm): 44

Remaining Lease (in years): 58

Closest Hawker Distance: 0.19

Number of Hawkers within 10m: 0

Closest MRT Distance: 0.8988

Number of MRTs within 10m: 4

Closest Park Distance: 0.874

Number of Parks within 10m: 23

Closest School Distance: 0.331

Number of Schools within 10m: 20

Closest Shopping Mall Distance: 1.001

Number of Shopping Malls within 10m: 8

Closest Supermarket Distance: 0.742

Number of Supermarkets within 10m: 10

Predict

Random Forest Prediction

Model Name: [1] 418102.3

Range: [1] 448102.5 451028.7

XGBoost Prediction

Model Name: [1] 377725.9

Range: [1] 355861.3 396248.5

Figure 4.10: Predicting housing resale price

Within the Shiny application, users have the flexibility to modify parameters via the interface on the left side. Upon selecting their desired criteria and clicking the "Predict" option, the application promptly displays predictions from both the Random Forest and XGBoost models. Additionally, to offer a broader perspective on the predicted housing resale prices, the application provides the interquartile range, showcasing the Q1 to Q3 range of the predictions, as illustrated in Figure 4.10.

5.0 Discussion

The analysis presented in the results section underscores the nuanced and multifaceted nature of the housing landscape in Singapore. This discussion seeks to delve deeper into these findings, drawing connections, implications, and potential future directions.

5.1 Synthesizing Model Outcomes

The performance metrics, chiefly the RMSE and MAPE, have validated the selection of Random Forest and XGBoost as pivotal models for this analysis. Their superiority in predicting resale prices can be attributed to their inherent capabilities in handling non-linearities and intricate interactions between variables.

However, beyond mere prediction, the feature importance derived from both models offers strategic insights. The consistent emphasis on features like 'floor area' and 'flat type' across both models underscores their universal significance in influencing resale prices. This congruence across models reinforces the reliability and validity of these findings.

5.2 Spatial Distribution and Urban Dynamics

The preliminary analysis highlighted stark contrasts in amenity accessibility across towns. Such disparities might be rooted in historical urban planning decisions, demographic shifts, or evolving residential preferences. Understanding these spatial dynamics is pivotal for urban planners and policymakers, aiding them in future infrastructure developments or resource allocations.

5.3 Integrating Analysis with User Experience

The Shiny application, as presented, epitomizes the amalgamation of rigorous data analysis with user-centric design. While the unsupervised recommendation system offers users

personalized housing suggestions, the supervised predictive model empowers them with resale price estimations. This dual approach not only facilitates informed decision-making but also democratizes data insights, making them accessible and actionable for the general populace.

5.4 Implications and Recommendations

The emphasis on 'floor area', 'flat type' suggests potential areas of focus for real estate developers or investors. For instance, properties with larger floor areas in desirable towns could command premium prices.

The Shiny application, in its current iteration, offers substantial utility to users. However, incorporating real-time data updates or expanding the range of amenities could enhance its predictive accuracy and recommendation relevance.

5.5 Limitations and Future Directions

While the models and analyses are robust, they are not devoid of limitations. Factors like historical price trends, imminent urban development projects, or socio-economic shifts were not incorporated into the current models but could influence resale prices.

Furthermore, evolving technologies, such as deep learning or neural networks, could be explored in future iterations to enhance predictive capabilities.

6.0 Conclusion

Characterized by a unique blend of spatial distribution and urban dynamics, Singapore's housing landscape offers a compelling field of study. Significant insights into the factors that influence the resale prices of housing units in the city-state have been gained through rigorous data analysis and modeling.

The adeptness of certain models, particularly Random Forest and XGBoost, in predicting resale prices has been proven. Not only have these models showcased impressive predictive capabilities, but key features, such as 'floor area' and 'flat type', have also been illuminated as being instrumental in determining resale prices.

Understanding has been further enriched by geospatial analysis, revealing the intricate balance of amenities across towns and emphasizing the importance of spatial context in urban studies.

A comprehensive platform to make informed housing decisions has been provided to users by offering both a recommendation system and a predictive tool.

It's crucial to acknowledge that, while many aspects of Singapore's housing market have been illuminated by our research, the urban fabric is ever-evolving. The factors influencing its housing market will continue to grow and change as the city does. For now, however, a significant step forward in understanding the current dynamics at play has been taken.

Looking ahead, immense potential exists to further refine our models, incorporate more variables, and leverage emerging technologies. Commitment to harnessing data for the benefit of its residents remains as the journey in understanding Singapore's urban landscape continues.

7.0 Recommendations

In light of the comprehensive analysis conducted and the tools developed, a structured approach is recommended for users and stakeholders to derive the most benefit. A detailed guide is provided here to ensure the utility of our findings and applications is maximized:

7.1 Structured Housing Exploration

7.1.1 Initial Exploration with Recommendation System

Begin your housing search by utilizing the recommendation system. This system, informed by historical data, offers a curated list of properties that align with user preferences. By starting here, users can gain insights into properties that have historically matched similar criteria and understand the sales patterns of such properties.

7.1.2 Informed Decision Making through the Resale Price Prediction Tool

After narrowing down potential housing options, transition to the resale price prediction tool. This tool provides a granular understanding of potential transaction values based on various influential factors. By leveraging this, users can make informed decisions, ensuring they get optimal value for their investments.

7.2 Continuous Data Updates and Training

The housing market is inherently fluid, with prices shaped by an ever-evolving array of factors. Given this dynamism, it is imperative to maintain an up-to-date dataset. By consistently integrating the most recent quarterly housing resale prices and retraining the models on this fresh data, predictions can be fine-tuned to anticipate the upcoming quarter's price trends more accurately. Stakeholders and developers of the application must recognize

the significance of these periodic updates to ensure that the tool remains both relevant and reflective of the prevailing market conditions.

7.3 Diversify Data Sources

To further enhance the accuracy and relevance of predictions, consider incorporating additional data sources. Integrating macroeconomic indicators, upcoming infrastructural projects, or even socio-demographic trends can offer a more holistic view of the housing market and refine predictions.

In conclusion, the tools and insights provided aim to empower users in their housing decisions. By adhering to these recommendations, users can ensure they navigate Singapore's housing market with clarity, confidence, and a data-driven edge.

List of References

- [1] T. Lau, "Singapore HDB Flat Resale Prices (1990-2020)," Kaggle, 2020. [Online]. Available: <https://www.kaggle.com/datasets/teyang/singapore-hdb-flat-resale-prices-19902020>
- [2] PropertyGuru. [Online]. Available: <https://www.propertyguru.com.sg/>
- [3] T. Vincenty, "Direct and inverse solutions of geodesics on the ellipsoid with application of nested equations," *Survey Review*, vol. 23, no. 176, pp. 88-93, 1975.
- [4] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016.
- [5] A. K. Jain, "Data clustering: 50 years beyond K-means," *Pattern Recognition Letters*, vol. 31, no. 8, pp. 651-666, 2010.
- [6] J. MacQueen, "Some methods for classification and analysis of multivariate observations," in *Proc. Fifth Berkeley Symp. on Math. Statist. and Prob.*, vol. 1, no. 14, pp. 281-297, 1967.
- [7] R. Sibson, "SLINK: an optimally efficient algorithm for the single-link cluster method," *The Computer Journal*, vol. 16, no. 1, pp. 30-34, 1973.
- [8] M. Ester, H. P. Kriegel, J. Sander, and X. Xu, "A density-based algorithm for discovering clusters in large spatial databases with noise," in *Proc. Kdd*, vol. 96, pp. 226-231, 1996.
- [9] C. M. Bishop, *Pattern Recognition and Machine Learning*. Springer, 2006.
- [10] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5-32, 2001.

- [11] T. Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," in Proc. 22nd ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining, pp. 785-794, 2016.
- [12] C. Cortes and V. Vapnik, "Support-vector networks," Machine Learning, vol. 20, no. 3, pp. 273-297, 1995.
- [13] J. H. Friedman, "Greedy function approximation: a gradient boosting machine," Annals of Statistics, pp. 1189-1232, 2001.
- [14] R. Tibshirani, "Regression shrinkage and selection via the lasso," Journal of the Royal Statistical Society: Series B (Methodological), vol. 58, no. 1, pp. 267-288, 1996.
- [15] A. E. Hoerl and R. W. Kennard, "Ridge regression: Biased estimation for nonorthogonal problems," Technometrics, vol. 12, no. 1, pp. 55-67, 1970.
- [16] R. J. Hyndman and A. B. Koehler, "Another look at measures of forecast accuracy," International Journal of Forecasting, vol. 22, no. 4, pp. 679-688, 2006.
- [17] S. Makridakis, "Accuracy measures: theoretical and practical concerns," International Journal of Forecasting, vol. 9, no. 4, pp. 527-529, 1993.
- [18] Data.Gov.sg, "Housing Resale Data," 2023. [Online]. Available: https://beta.data.gov.sg/datasets/189/resources/d_8b84c4ee58e3cfc0ece0d773c8ca6abc/view
- [19] Data.Gov.sg, "School Data," 2022. [Online]. Available: <https://beta.data.gov.sg/datasets/457/view>
- [20] Data.Gov.sg, "Hawker Centers," 2022. [Online]. Available: <https://beta.data.gov.sg/datasets/1445/view>

- [21] Hxchua, "Train Stations in Singapore," Data.world, 2020. [Online]. Available: <https://data.world/hxchua/train-stations-in-singapore>
- [22] Data.Gov.sg, "Supermarkets," 2020. [Online]. Available: <https://beta.data.gov.sg/datasets/1565/view>
- [23] K. Gangula, "Shopping Mall Coordinates," Kaggle, 2023. [Online]. Available: <https://www.kaggle.com/datasets/karthikgangula/shopping-mall-coordinates>
- [24] Data.Gov.sg, "Parks," 2022. [Online]. Available: https://beta.data.gov.sg/datasets/1491/resources/d_0542d48f0991541706b58059381a6eca/view
- [25] Ministry of Education, Singapore, "Primary One Registration: Distance," 2023. [Online]. Available: <https://www.moe.gov.sg/primary/p1-registration/distance>

Appendix A : Combining housing dataset with the 6 amenities

Figure A.1 illustrates the process of acquiring longitudes and latitudes for the original housing dataset. Subsequently, Figures A.2, A.6, and A.9 depict the codes used to obtain longitudes and latitudes for their respective amenities datasets. Finally, the rest of the figures present the codes implemented to calculate the number of amenities and identify the closest ones.

```
18 ~~~{r adding longitude and latitude to original q2 2023}
19 library(ggmap)
20
21 # Set your Google Maps API key
22 register_google(key = "AIzaSyCmIgb23RkKNuz79bEBPVtodeM17AK60E0")
23
24 # Read addresses from the CSV file
25 address_data <- read.csv("data2023.csv")
26
27 # Convert the 'month' column to character
28 address_data$month <- as.character(address_data$month)
29
30 # Filter data for months 4, 5, and 6
31 filtered_data <- address_data[address_data$month >= "2023-04" & address_data$month <= "2023-06", ]
32
33 # Initialize empty vectors to store latitudes and longitudes
34 latitudes <- numeric()
35 longitudes <- numeric()
36
37 # Iterate through each row and geocode the address
38 for (i in seq_len(nrow(filtered_data))) {
39   block <- filtered_data$block[i]
40   street <- filtered_data$street[i]
41   address <- paste(block, street)
42
43   result <- geocode(address)
44   if (!is.na(result$lat) && !is.na(result$lon)) {
45     latitudes <- c(latitudes, result$lat)
46     longitudes <- c(longitudes, result$lon)
47   } else {
48     latitudes <- c(latitudes, NA)
49     longitudes <- c(longitudes, NA)
50   }
51 }
52
53 # Add latitudes and longitudes to the filtered data
54 filtered_data$Latitude <- latitudes
55 filtered_data$Longitude <- longitudes
56
57 # Save the filtered data with latitudes and longitudes
58 write.csv(filtered_data, file = "filtered_geocoded_data.csv", row.names = FALSE)
59
60 cat("Geocoding and filtering completed. Filtered data saved to filtered_geocoded_data.csv.\n")
61
62 ~~~
63
64
```

52:2 Chunk 2: adding longitude and latitude to original q2 2023 R Markdown

Figure A.1: Adding longitude and latitude to original Q2 2023

```

66
67 - ```{r getting parks}
68
69 library(sf)
70
71 # Read KML data
72 kml_data <- st_read("Parks.kml")
73
74 # Extract coordinates
75 coordinates <- st_coordinates(kml_data)
76
77 # Optional: Save coordinates to a CSV file
78 write.csv(coordinates, file = "parks.csv", row.names = FALSE)
79
80 cat("Coordinates extracted and saved to coordinates.csv.\n")
81
82 ```
83
84
85
86 - ```{r getting hawker}
87 library(sf)
88
89 # Read KML data
90 kml_data <- st_read("HawkerCentresKML.kml")
91
92 # Extract coordinates
93 coordinates <- st_coordinates(kml_data)
94
95 # Optional: Save coordinates to a CSV file
96 write.csv(coordinates, file = "Hawker.csv", row.names = FALSE)
97
98 cat("Coordinates extracted and saved to coordinates.csv.\n")
99
100 ```
101

```

Figure A.2: Getting parks and hawkers longitude and latitudes

```

L01
L02 - ```{r adding hawker}
L03 # Load necessary libraries
L04 library(dplyr)
L05 library(geosphere)
L06
L07 # Read the CSV data
L08 data <- read.csv("filtered_geocoded_data.csv") # Replace 'your_data.csv' with your actual data file name
L09 hawker_data <- read.csv("Hawker.csv") # Replace 'hawker_data.csv' with the Hawker data file name
L10
L11 # Remove rows with negative longitude and latitude
L12 data <- data %>%
L13   filter(Longitude >= 0, Latitude >= 0)
L14
L15 # Extract Hawker coordinates
L16 hawker_coors <- hawker_data[, c("X", "Y")]
L17
L18 # Create a new column for closest distance
L19 data <- data %>%
L20   mutate(hawker_closest_distance = 0.0,
L21          num_hawker_1km = 0)
L22
L23 - for (i in 1:nrow(data)) {
L24   flat_coors <- c(data[i, "Longitude"], data[i, "Latitude"])
L25   distances <- distVincentysphere(flat_coors, hawker_coors)
L26   min_distance <- min(distances) / 1000 # Convert meters to kilometers
L27   data[i, "hawker_closest_distance"] <- min_distance
L28
L29   # Count Hawker locations within 1km
L30   num_hawker_within_1km <- sum(distances <= 1000) # Check distances <= 1000 meters (1km)
L31   data[i, "num_hawker_1km"] <- num_hawker_within_1km
L32 }
L33
L34 # Save the updated data to a new CSV file
L35 write.csv(data, "updated_data.csv", row.names = FALSE)
L36

```

Figure A.3: Adding number of hawkers within 1km radius , and the closest hawker distance

```

137/
138- ```{r adding mrt}
139 # Load necessary libraries
140 library(dplyr)
141 library(geosphere)
142
143 # Read the CSV data
144 data <- read.csv("updated_data.csv") # Replace with the actual updated data file name
145 mrt_data <- read.csv("mrt.csv") # Replace with the MRT data file name
146
147 # Extract MRT coordinates
148 mrt_coords <- mrt_data[, c("Longitude", "Latitude")]
149
150 # Create new columns for closest MRT distance and number of MRT locations within 2km
151 data <- data %>%
152   mutate(mrt_closest_dist = 0.0,
153          num_mrt_1km = 0) # Initialize the new columns
154
155 # Calculate closest MRT distance and count MRT locations within 1km
156 for (i in 1:nrow(data)) {
157   flat_coords <- c(data[i, "Longitude"], data[i, "Latitude"])
158   distances <- distvincentysphere(flat_coords, mrt_coords)
159   min_distance <- min(distances) / 1000 # Convert meters to kilometers
160   data[i, "mrt_closest_dist"] <- min_distance
161
162   # Count MRT locations within 2km
163   num_mrt_within_1km <- sum(distances <= 1000) # Check distances <= 1000 meters (1km)
164   data[i, "num_mrt_1km"] <- num_mrt_within_1km
165 }
166
167 # Save the updated data to a new CSV file
168 write.csv(data, "updated_data.csv", row.names = FALSE)
169
170- ```

```

Figure A.4: Adding number of mrt within 1km radius , and the closest mrt distance

```

172- ```{r adding park}
173 # Load necessary libraries
174 library(dplyr)
175 library(geosphere)
176
177 # Read the CSV data
178 data <- read.csv("updated_data.csv") # Replace with the actual updated data file name
179 park_data <- read.csv("park.csv") # Replace with the park data file name
180
181 # Extract park coordinates
182 park_coords <- park_data[, c("X", "Y")]
183
184 # Create new columns for closest park distance and number of park locations within 2km
185 data <- data %>%
186   mutate(park_closest_dist = 0.0,
187          num_park_1km = 0) # Initialize the new columns
188
189 # Calculate closest park distance and count park locations within 1km
190 for (i in 1:nrow(data)) {
191   flat_coords <- c(data[i, "Longitude"], data[i, "Latitude"])
192   distances <- distvincentysphere(flat_coords, park_coords)
193   min_distance <- min(distances) / 1000 # Convert meters to kilometers
194   data[i, "park_closest_dist"] <- min_distance
195
196   # Count park locations within 2km
197   num_park_within_1km <- sum(distances <= 1000) # Check distances <= 1000 meters (1km)
198   data[i, "num_park_1km"] <- num_park_within_1km
199 }
200
201 # Save the updated data to a new CSV file
202 write.csv(data, "updated_data.csv", row.names = FALSE)
203
204- ```
205-

```

Figure A.5: Adding number of park within 1km radius , and the closest park distance

```

206 ~~~{r adding coordinates to school}
207 # Load necessary libraries
208 library(ggmap)
209
210 # Set your Google Maps API key
211 register_google(key = "AIzaSyCmIgB23RKKNuz79bEBPVtodeM17AK60E0")
212
213 # Read the CSV data
214 school_data <- read.csv("school.csv")
215
216 # Define a function to geocode postal codes to latitude and longitude using ggmap
217 ~ geocode_postal <- function(postal_code) {
218   geo_data <- geocode(postal_code)
219   return(geo_data)
220 }
221
222 # Create new columns for latitude and longitude
223 school_data$Latitude <- NA
224 school_data$Longitude <- NA
225
226 # Iterate through each row and geocode the postal code
227 ~ for (i in 1:nrow(school_data)) {
228   postal_code <- as.character(school_data[i, "postal_code"])
229   ~ if (!is.na(postal_code)) {
230     geo_data <- geocode_postal(postal_code)
231     school_data[i, "Latitude"] <- geo_data$lat
232     school_data[i, "Longitude"] <- geo_data$lon
233   }
234 }
235
236 # Save the geocoded data to a new CSV file
237 write.csv(school_data, "school_with_coords.csv", row.names = FALSE)
238
239 cat("Geocoding completed. Geocoded data saved to school_with_coords.csv.\n")
240
241 ~~~
242

```

Figure A.6: Adding longitude and latitude in the school dataset from the postal code

```

243 ~~~{r adding school}
244 # Load necessary libraries
245 library(dplyr)
246 library(geosphere)
247
248 # Read the CSV data
249 data <- read.csv("updated_data.csv") # Replace with the actual updated data file name
250 school_data <- read.csv("school_with_coords.csv") # Replace with the school data file name
251
252 # Extract school coordinates
253 school_coords <- na.omit(school_data[, c("Longitude", "Latitude")])
254
255 # Create new columns for closest school distance and number of school locations within 2km
256 data <- data %>%
257   mutate(sch_closest_dist = 0.0,
258          num_sch_2km = 0) # Initialize the new columns
259
260 # Calculate closest school distance and count school locations within 2km
261 ~ for (i in 1:nrow(data)) {
262   flat_coords <- c(data[i, "Longitude"], data[i, "Latitude"])
263   distances <- distVincentysphere(flat_coords, school_coords)
264   min_distance <- min(distances) / 1000 # Convert meters to kilometers
265   data[i, "sch_closest_dist"] <- min_distance
266
267   # Count school locations within 2km
268   num_sch_within_2km <- sum(distances <= 2000) # Check distances <= 2000 meters (2km)
269   data[i, "num_sch_2km"] <- num_sch_within_2km
270 }
271
272 # Save the updated data to a new CSV file
273 write.csv(data, "updated_data.csv", row.names = FALSE)
274
275 ~~~
276

```

Figure A.7: Adding number of school within 2km radius , and the closest school distance


```

276 |
277 | ```{r adding shoppingmall}
278 | # Load necessary libraries
279 | library(dplyr)
280 | library(geosphere)
281 |
282 | # Read the CSV data
283 | data <- read.csv("updated_data.csv") # Replace with the actual updated data file name
284 | shoppingmall_data <- read.csv("shoppingmall.csv") # Replace with the shopping mall data file name
285 |
286 | # Extract shopping mall coordinates and drop rows with NA values
287 | shoppingmall_coors <- na.omit(shoppingmall_data[, c("LONGITUDE", "LATITUDE")])
288 |
289 | # Create new columns for closest shopping mall distance and number of shopping malls within 2km
290 | data <- data %>%
291 |   mutate(shoppingmall_closest_dist = 0.0,
292 |          num_shoppingmall_1km = 0) # Initialize the new columns
293 |
294 | # Calculate closest shopping mall distance and count shopping malls within 1km
295 | for (i in 1:nrow(data)) {
296 |   flat_coors <- c(data[i, "Longitude"], data[i, "Latitude"])
297 |   distances <- distVincentySphere(flat_coors, shoppingmall_coors)
298 |   min_distance <- min(distances) / 1000 # Convert meters to kilometers
299 |   data[i, "shoppingmall_closest_dist"] <- min_distance
300 |
301 |   # Count shopping malls within 2km
302 |   num_shoppingmall_within_1km <- sum(distances <= 1000) # Check distances <= 1000 meters (2km)
303 |   data[i, "num_shoppingmall_1km"] <- num_shoppingmall_within_1km
304 | }
305 |
306 | # Save the updated data to a new CSV file
307 | write.csv(data, "updated_data.csv", row.names = FALSE)
308 |
309 | cat("Processing completed. Updated data saved to updated_data.csv.\n")
310 |
311 | ```
312 |

```

Figure A.8: Adding number of shopping mall within 1km radius , and the closest shopping mall distance

```

313 |
314 | ```{r adding coordinates to supermarket}
315 | # Read the CSV data
316 | supermarket_data <- read.csv("supermarket.csv") # Replace with the actual supermarket data file name
317 |
318 | # Extract postal codes using regular expressions
319 | postal_codes <- gsub(".*S\\((\\d+\\)\\).*", "\\1", supermarket_data$premise_address)
320 |
321 | # Create a new column for postal codes in the supermarket data
322 | supermarket_data$PostalCode <- postal_codes
323 |
324 | # Set your Google Maps API key
325 | register_google(key = "AIzaSyCmIgB23RKNuz79bEBPVtodeM17AK6OE0")
326 |
327 | # Define a function to geocode postal codes to latitude and longitude using ggmap
328 | geocode_postal <- function(postal_code) {
329 |   geo_data <- geocode(postal_code)
330 |   return(geo_data)
331 | }
332 |
333 | # Create new columns for latitude and longitude
334 | supermarket_data$Latitude <- NA
335 | supermarket_data$Longitude <- NA
336 |
337 | # Iterate through each row and geocode the postal code
338 | for (i in 1:nrow(supermarket_data)) {
339 |   postal_code <- as.character(supermarket_data[i, "PostalCode"])
340 |   if (!is.na(postal_code)) {
341 |     geo_data <- geocode_postal(postal_code)
342 |     supermarket_data[i, "Latitude"] <- geo_data$lat
343 |     supermarket_data[i, "Longitude"] <- geo_data$lon
344 |   }
345 | }
346 |
347 | # Save the updated supermarket data to a new CSV file
348 | write.csv(supermarket_data, "supermarket_with_coors.csv", row.names = FALSE)
349 |
350 | cat("Processing completed. Supermarket data with postal codes saved to supermarket_with_postal.csv.\n")
351 |
352 | ```
353 |

```

Figure A.9: Adding longitude and latitude in the supermarket dataset from the postal code

```

353
354 - ```{r adding supermarket}
355 # Load necessary libraries
356 library(dplyr)
357 library(geosphere)
358
359 # Read the CSV data
360 data <- read.csv("updated_data.csv") # Replace with the actual updated data file name
361 supermarket_data <- read.csv("supermarket_with_coords.csv") # Replace with the supermarket data file name
362
363 # Extract supermarket coordinates and drop rows with NA values
364 supermarket_coords <- na.omit(supermarket_data[, c("Longitude", "Latitude")])
365
366 # Create new columns for closest supermarket distance and number of supermarkets within 2km
367 data <- data %>%
368   mutate(supermarket_closest_dist = 0.0,
369          num_supermarket_1km = 0) # Initialize the new columns
370
371 # Calculate closest supermarket distance and count supermarkets within 2km
372 - for (i in 1:nrow(data)) {
373   flat_coords <- c(data[i, "Longitude"], data[i, "Latitude"])
374   distances <- distvincentysphere(flat_coords, supermarket_coords)
375   min_distance <- min(distances) / 1000 # Convert meters to kilometers
376   data[i, "supermarket_closest_dist"] <- min_distance
377
378   # Count supermarkets within 1km
379   num_supermarket_within_1km <- sum(distances <= 1000) # Check distances <= 1000 meters (1km)
380   data[i, "num_supermarket_1km"] <- num_supermarket_within_1km
381 - }
382
383 # Save the updated data to a new CSV file
384 write.csv(data, "updated_data.csv", row.names = FALSE)
385
386 cat("Processing completed. Updated data saved to updated_data.csv.\n")
387
388 - ```
389
390

```

Figure A.10: Adding number of supermarket within 1km radius , and closest supermarket distance

Appendix B: Preliminary analysis of the combined dataset

Figure B.1 pertains to the data cleaning process, focusing on the elimination of outliers and the selection of meaningful columns for subsequent use. Figures B.2 and B.3 are dedicated to performing statistical analysis and generating visual plots, respectively, analyzing the housing dataset grouped by town.

```
11
12- ```{r data cleaning}
13 # Read data
14 library(tidyverse)
15 data <- read_csv('updated_data.csv', show_col_types = FALSE)
16 # Data cleaning:
17 # 1. Drop specified columns
18 # 2. Remove rows where flat_model is "2-room"
19 # 3. Remove rows where distances are greater than 4km for all towns
20 cleaned_data <- data %>%
21   select(-block, -street_name, -month, -lease_commence_date, -Latitude, -Longitude) %>%
22   filter(flat_model != "2-room") %>%
23   filter(hawker_closest_distance <= 4 &
24     mrt_closest_dist <= 4 &
25     park_closest_dist <= 4 &
26     sch_closest_dist <= 4 &
27     shoppingmall_closest_dist <= 4 &
28     supermarket_closest_dist <= 4)
29
30 # Save the cleaned data to a new CSV file
31 write_csv(cleaned_data, "final.csv", row.names = FALSE)
32- ```
33
```

Figure B.1: Data cleaning to include the relevant columns for future data analysis

```

34
35 # [r preliminary analysis]
36 library(dplyr)
37 library(tidyr)
38 library(ggplot2)
39 library(grid)
40 library(gridExtra)
41
42 # Read data
43 data <- read.csv("final.csv")
44
45 # Create a function to obtain the summary statistics
46 summary_stats <- function(column) {
47   c(Min = min(column, na.rm = TRUE),
48     Median = median(column, na.rm = TRUE),
49     Mean = mean(column, na.rm = TRUE),
50     Max = max(column, na.rm = TRUE),
51     '1st Qu.' = quantile(column, 0.25, na.rm = TRUE),
52     '3rd Qu.' = quantile(column, 0.75, na.rm = TRUE))
53 }
54
55 # Apply the function to each column of interest
56 summaries <- data.frame(
57   t(sapply(data[, c("hawker_closest_distance",
58     "num_hawker_1km",
59     "num_closest_dist",
60     "num_mrt_1km",
61     "park_closest_dist",
62     "num_park_1km",
63     "sch_closest_dist",
64     "num_sch_1km",
65     "shoppingmall_closest_dist",
66     "num_shoppingmall_1km",
67     "supermarket_closest_dist",
68     "num_supermarket_1km")], summary_stats))
69 )
70
71 # View the summary table
72 print(summaries)
73
74 # Convert the data frame to a grid table
75 grid_table <- tableGrob(summaries)
76
77 # Save the table as an image
78 ggsave(filename = "summary_table.png", plot = grid_table, width = 20, height = 8)
79
80

```

Figure B.2: Preliminary analysis part 1

```

81
82 # calculate the average distances by town in meters for all amenities
83 avg_dist_by_town <- data %>%
84   group_by(town) %>%
85   summarise(
86     avg_dist_hawker = mean(hawker_closest_distance, na.rm = TRUE),
87     avg_dist_mrt = mean(mrt_closest_dist, na.rm = TRUE),
88     avg_dist_park = mean(park_closest_dist, na.rm = TRUE),
89     avg_dist_sch = mean(sch_closest_dist, na.rm = TRUE),
90     avg_dist_shoppingmall = mean(shoppingmall_closest_dist, na.rm = TRUE),
91     avg_dist_supermarket = mean(supermarket_closest_dist, na.rm = TRUE)
92   )
93
94 # convert data to long format for plotting
95 avg_dist_long <- avg_dist_by_town %>%
96   pivot_longer(cols = starts_with("avg_dist"),
97     names_to = "amenity",
98     values_to = "avg_distance") %>%
99   mutate(amenity = recode(amenity,
100     avg_dist_hawker = "hawker center",
101     avg_dist_mrt = "mrt station",
102     avg_dist_park = "park",
103     avg_dist_sch = "school",
104     avg_dist_shoppingmall = "shopping mall",
105     avg_dist_supermarket = "supermarket"))
106
107 g <- ggplot(avg_dist_long, aes(x = town, y = avg_distance, fill = amenity)) +
108   geom_bar(stat = "identity", position = position_dodge(width = 0.8), width = 0.7) + # Adjust bar width and dodge width
109   facet(
110     title = "Average distance to amenities by town",
111     y = "Average distance (in km)",
112     x = "town",
113     fill = "amenity"
114   )
115 coord_flip() +
116   theme(text = element_text(angle = 0))
117 ggsave("avg_distance_plot_spaced_amenities.png", plot = g, width = 20, height = 22) # Adjusted height to accommodate spacing
118
119 # calculate the average number of facilities by town for all amenities
120 avg_num_by_town <- data %>%
121   group_by(town) %>%
122   summarise(
123     avg_num_hawker = mean(num_hawker_1km, na.rm = TRUE),
124     avg_num_mrt = mean(num_mrt_1km, na.rm = TRUE),
125     avg_num_park = mean(num_park_1km, na.rm = TRUE),
126     avg_num_sch = mean(num_sch_1km, na.rm = TRUE), # note that schools use a 2km radius
127     avg_num_shoppingmall = mean(num_shoppingmall_1km, na.rm = TRUE),
128     avg_num_supermarket = mean(num_supermarket_1km, na.rm = TRUE)
129   )
130
131 # convert data to long format for plotting
132 avg_num_long <- avg_num_by_town %>%
133   pivot_longer(cols = starts_with("avg_num"),
134     names_to = "amenity",
135     values_to = "avg_number") %>%
136   mutate(amenity = recode(amenity,
137     avg_num_hawker = "hawker center",
138     avg_num_mrt = "mrt station",
139     avg_num_park = "park",
140     avg_num_sch = "school",
141     avg_num_shoppingmall = "shopping mall",
142     avg_num_supermarket = "supermarket"))
143
144 # create the plot with increased space between bars
145 g <- ggplot(avg_num_long, aes(x = town, y = avg_number, fill = amenity)) +
146   geom_bar(stat = "identity", position = position_dodge(width = 0.8), width = 0.7) +
147   facet(
148     title = "Average number of facilities by town",
149     y = "Average number of facilities",
150     x = "town",
151     fill = "amenity"
152   )
153 coord_flip() +
154   theme(text = element_text(angle = 0))
155 ggsave("avg_number_plot_spaced_amenities.png", plot = g, width = 20, height = 22) # Adjusted height for spacing
156
157
158
159

```

Figure B.3: Preliminary analysis part 2

Appendix C: unsupervised and supervised analysis of the combined dataset

Figures C.1 to C.6 depict the code for unsupervised machine learning in the shiny app, while Figures C.7 to C.11 represent supervised machine learning code for the app.

```
160
161- '''{r with map }
162 library(shiny)
163 library(caret) # For one-hot encoding
164 library(tidyverse)
165 library(fpc)
166 library(leaflet) # Added For map integration
167
168 # Read data
169 data <- read.csv("updated_data.csv")
170
171 # Extract years
172 data$years <- as.numeric(gsub(" years.", "", data$remaining_lease))
173
174 # Extract months
175 data$months <- 0
176 data$months[grepl("month", data$remaining_lease)] <- as.numeric(gsub(".", "", gsub(" month.", "", data$remaining_lease[grepl("month", data$remaining_lease)])))
177
178
179 ui <- fluidPage(
180   titlePanel("Housing Recommendations"),
181
182   # Row for Selection and Map
183   fluidRow(
184     # Left Panel for Selection
185     column(2,
186       checkboxGroupInput("selected_parameters", "Select Parameters:",
187         choices = c("MRT", "Park", "School", "Supermarket", "Shopping Mall", "Hawker Center", "Town", "Flat Type", "Storey Range", "Flat Model", "Floor
188         Area", "Remaining Lease"),
189         selected = "MRT"),
190
191       radioButtons("clustering_method", "Select Clustering Method:",
192         choices = c("K-means Clustering", "Hierarchical Clustering", "DBSCAN"),
193         selected = "K-means Clustering"),
194
195       numericInput("k_clusters", "Number of Clusters:", value = 2, min = 1, max = 10000, step = 1),
196
197       conditionalPanel(
198         condition = "input.clustering_method == 'DBSCAN'",
199         numericInput("dbscan_eps", "EPS (Max Distance):", value = 0.5, min = 0, max = 10, step = 0.1),
200         numericInput("dbscan_minPts", "Min Points:", value = 5, min = 1, max = 100, step = 1)
201       ),
202
203       # Numeric inputs for distances
204       conditionalPanel(
205         condition = "'MRT' %in% input.selected_parameters",
206         numericInput("mrt_distance", "At most MRT Distance (km):", value = 1, min = 0, max = 10, step = 0.1),
207         numericInput("min_run_mrts", "Minimum Number of MRT Stations (1Km) radius:", value = 1, min = 0, max = 10, step = 1)
208       ),
209       conditionalPanel(
210         condition = "'Park' %in% input.selected_parameters",
211         numericInput("park_distance", "At most Park Distance (km):", value = 1, min = 0, max = 10, step = 0.1),
212         numericInput("min_run_parks", "Minimum Number of Parks (1Km) radius:", value = 1, min = 0, max = 10, step = 1)
213       ),
214       conditionalPanel(
215         condition = "'School' %in% input.selected_parameters",
216         numericInput("school_distance", "At most School Distance (km):", value = 1, min = 0, max = 10, step = 0.1),
217         numericInput("min_run_schools", "Minimum Number of Schools (2Km) radius:", value = 1, min = 0, max = 10, step = 1)
218       ),
219       conditionalPanel(
220         condition = "'Supermarket' %in% input.selected_parameters",
221         numericInput("supermarket_distance", "At most Supermarket Distance (km):", value = 1, min = 0, max = 10, step = 0.1),
222         numericInput("min_run_supermarkets", "Minimum Number of Supermarkets (1Km) radius:", value = 1, min = 0, max = 10, step = 1)
223       )
224     )
225   )
226 )
```

Figure C.1: unsupervised recommendation system shiny app part 1

```

222 conditionalPanel(
223   condition = "'Shopping Mall' %in% input.selected_parameters",
224   numericInput("shoppingmall_distance", "At most Shopping Mall Distance (km):", value = 1, min = 0, max = 10, step = 0.1),
225   numericInput("min_num_shoppingmalls", "Minimum Number of Shopping Malls (1Km) radius:", value = 1, min = 0, max = 10, step = 1)
226 ),
227 conditionalPanel(
228   condition = "'Hawker Center' %in% input.selected_parameters",
229   numericInput("hawker_distance", "At most Hawker Center Distance (km):", value = 1, min = 0, max = 10, step = 0.1),
230   numericInput("min_num_hawkers", "Minimum Number of Hawker Centers (1Km) radius:", value = 1, min = 0, max = 10, step = 1)
231 ),
232
233 # Price filter
234 checkboxInput("apply_price_filter", "Apply Price Filter", value = FALSE),
235 conditionalPanel(
236   condition = "input.apply_price_filter == true",
237   numericInput("min_price", "Minimum Price:", value = 200000, min = 0, max = 2000000, step = 10000),
238   numericInput("max_price", "Maximum Price:", value = 1000000, min = 0, max = 2000000, step = 10000)
239 ),
240
241 # Conditional input for floor_area_sqm based on checkbox
242 conditionalPanel(
243   condition = "'Floor Area' %in% input.selected_parameters",
244   sliderInput("floor_area_sqm_range", "Floor Area (sqm):", min = min(data$floor_area_sqm), max = max(data$floor_area_sqm), value = c(min(data$floor_area_sqm),
245     max(data$floor_area_sqm)))
246 ),
247
248 # Conditional inputs for remaining_lease based on checkbox
249 conditionalPanel(
250   condition = "'Remaining Lease' %in% input.selected_parameters",
251   numericInput("remaining_lease_years", "Remaining Lease (years):", value = max(data$years), min = min(data$years), max = max(data$years)),
252   numericInput("remaining_lease_months", "Remaining Lease (months):", value = 0, min = 0, max = 11)
253 ),
254
255 # Dropdowns for categorical features
256 conditionalPanel(
257   condition = "'Town' %in% input.selected_parameters",
258   selectInput("town", "Select Town:", choices = unique(data$town))
259 ),
260 conditionalPanel(
261   condition = "'Flat Type' %in% input.selected_parameters",
262   selectInput("flat_type", "Select Flat Type:", choices = unique(data$flat_type))
263 ),
264 conditionalPanel(
265   condition = "'Storey Range' %in% input.selected_parameters",
266   selectInput("storey_range", "Select Storey Range:", choices = unique(data$storey_range))
267 ),
268 conditionalPanel(
269   condition = "'Flat Model' %in% input.selected_parameters",
270   selectInput("flat_model", "Select Flat Model:", choices = unique(data$flat_model))
271 ),
272
273 condition = "'MRT' %in% input.selected_parameters || 'Park' %in% input.selected_parameters || 'School' %in% input.selected_parameters || 'Supermarket' %in%
274 input.selected_parameters || 'Shopping Mall' %in% input.selected_parameters || 'Hawker Center' %in% input.selected_parameters",
275 selectInput("numeric_criterion", "Select Numeric Criterion for Recommendation:",
276   choices = c("At most MRT Distance" = "mrt_closest_dist",
277     "Minimum Number of MRT Stations" = "num_mrt_1km",
278     "At most Park Distance" = "park_closest_dist",
279     "Minimum Number of Parks" = "num_park_1km",
280     "At most School Distance" = "sch_closest_dist",
281     "Minimum Number of Schools" = "num_sch_2km",
282     "At most Supermarket Distance" = "supermarket_closest_dist",

```

Figure C.2: unsupervised recommendation system shiny app part 2

```

277       "Minimum Number of Parks" = "num_park_1km",
278       "At most School Distance" = "sch_closest_dist",
279       "Minimum Number of Schools" = "num_sch_2km",
280       "At most Supermarket Distance" = "supermarket_closest_dist",
281       "Minimum Number of Supermarkets" = "num_supermarket_1km",
282       "At most Shopping Mall Distance" = "shoppingmall_closest_dist",
283       "Minimum Number of Shopping Malls" = "num_shoppingmall_1km",
284       "At most Hawker Center Distance" = "hawker_closest_distance",
285       "Minimum Number of Hawker Centers" = "num_hawker_1km",
286       "Minimum Mean Resale Price" = "resale_price_mean"),
287       selected = "mrt_closest_dist")
288     },
289   ),
290   # Right Panel for Map
291   column(10,
292     tableOutput("Filtered_data")
293   ),
294   # Below Panel for Recommendations
295   fluidRow(
296     column(12,
297       leafletOutput("map", height = "600px")
298     )
299   )
300 )
301 )
302 )
303
304 server <- function(input, output, session) {
305   recommendation_data <- reactive({
306     filtered_data <- data
307
308     # Filter based on selected parameters and their corresponding inputs
309     if ("MRT" %in% input$selected_parameters) {
310       filtered_data <- filtered_data[filtered_data$mrt_closest_dist <= input$mrt_distance, ]
311       if (nrow(filtered_data) == 0) return("No data after MRT filter")
312     }
313
314     if ("Park" %in% input$selected_parameters) {
315       filtered_data <- filtered_data[filtered_data$park_closest_dist <= input$park_distance, ]
316       if (nrow(filtered_data) == 0) return("No data after Park filter")
317     }
318
319     if ("School" %in% input$selected_parameters) {
320       filtered_data <- filtered_data[filtered_data$sch_closest_dist <= input$school_distance, ]
321       if (nrow(filtered_data) == 0) return("No data after School filter")
322     }
323
324     if ("Supermarket" %in% input$selected_parameters) {
325       filtered_data <- filtered_data[filtered_data$supermarket_closest_dist <= input$supermarket_distance, ]
326       if (nrow(filtered_data) == 0) return("No data after Supermarket filter")
327     }
328
329     if ("Shopping Mall" %in% input$selected_parameters) {
330       filtered_data <- filtered_data[filtered_data$shoppingmall_closest_dist <= input$shoppingmall_distance, ]
331       if (nrow(filtered_data) == 0) return("No data after Shopping Mall filter")
332     }
333
334     if ("Hawker Center" %in% input$selected_parameters) {
335       filtered_data <- filtered_data[filtered_data$hawker_closest_dist <= input$hawker_distance, ]
336       if (nrow(filtered_data) == 0) return("No data after Hawker Center filter")
337     }
338
339     # Filter by price range ONLY if "apply_price_filter" is TRUE
340     if (input$apply_price_filter) {

```

Figure C.3: unsupervised recommendation system shiny app part 3

```

401 if ("Shopping Mall" %in% input$selected_parameters && !is.null(input$min_num_shoppingmalls)) {
402   filtered_data <- filtered_data[filtered_data$num_shoppingmall_1km >= input$min_num_shoppingmalls, ]
403   if (nrow(filtered_data) == 0) return("No data after minimum shopping malls filter")
404 }
405
406 # Similarly add conditions for other amenities like schools
407 if ("School" %in% input$selected_parameters && !is.null(input$min_num_schools)) {
408   filtered_data <- filtered_data[filtered_data$num_sch_2km >= input$min_num_schools, ]
409   if (nrow(filtered_data) == 0) return("No data after minimum schools filter")
410 }
411
412 # ... Add conditions for additional amenities as needed
413
414 # Define the function to determine if the criterion is based on minimizing or maximizing
415 should_pick_min <- function(criterion) {
416   min_based_criteria <- c("mrt_closest_dist", "park_closest_dist", "sch_closest_dist",
417     "supermarket_closest_dist", "shoppingmall_closest_dist",
418     "hawker_closest_distance", "resale_price_mean")
419   return(criterion %in% min_based_criteria)
420 }
421
422 # Mapping between user input and the actual dataset variable name
423 user_to_data_map <- list("Town" = "town",
424   "Flat Type" = "flat_type",
425   "Storey Range" = "storey_range",
426   "Flat Model" = "flat_model")
427
428 selected_features <- c("town", "flat_type", "storey_range", "flat_model") # default features
429
430 # Check each feature: if it has only one unique value, exclude it from encoding
431 for (feature in selected_features) {
432   if (length(unique(filtered_data[[feature]])) == 1) {
433     selected_features <- setdiff(selected_features, feature)
434   }
435 }
436
437 # Convert the vector of features to a formula
438 formula <- as.formula(paste("~", paste(selected_features, collapse = " + "), "- 1"))
439
440 # Encode the features
441 encoded_features <- model.matrix(formula, filtered_data)
442
443 numeric_features <- filtered_data %>%
444   dplyr::select(
445     floor_area_sqm,
446     resale_price,
447     years,
448     num_hawker_1km,
449     mrt_closest_dist,
450     hawker_closest_distance,
451     shoppingmall_closest_dist,
452     park_closest_dist,
453     sch_closest_dist,
454     supermarket_closest_dist,
455     num_supermarket_1km,
456     num_mrt_1km,
457     num_sch_2km,
458     num_park_1km,
459     num_shoppingmall_1km,
460     months
461   )
462 scaled_numeric_features <- scale(numeric_features)
463 combined_features <- cbind(encoded_features, scaled_numeric_features)
464
465

```

Figure C.4: unsupervised recommendation system shiny app part 4


```

466 # Check for NA, NaN, or Inf values and remove corresponding rows
467 combined_features <- combined_features[!apply(is.na(combined_features) | is.infinite(combined_features), 1, any),]
468
469 # Apply the clustering methods
470 if(input$clustering_method == "K-means Clustering") {
471   set.seed(123)
472   clustering_model <- kmeans(combined_features, centers = input$k_clusters)
473
474   # Add additional logic for the mean resale price
475   if(input$numeric_criterion == "resale_price_mean") {
476     cluster_means <- tapply(filtered_data$resale_price, clustering_model$cluster, mean)
477   } else {
478     cluster_means <- clustering_model$centers[, input$numeric_criterion]
479   }
480
481   # Decide the recommended cluster
482   if (should_pick_min(input$numeric_criterion)) {
483     recommended_cluster <- which.min(cluster_means)
484   } else {
485     recommended_cluster <- which.max(cluster_means)
486   }
487
488   recommendation <- filtered_data[clustering_model$cluster == recommended_cluster, ]
489   # Limit the recommendation to the top 45 rows
490   recommendation <- head(recommendation, 45)
491 } else if(input$clustering_method == "Hierarchical Clustering") {
492   hc <- hclust(dist(combined_features))
493   cluster_assignments <- cutree(hc, k = input$k_clusters)
494
495   cluster_means <- tapply(filtered_data[input$numeric_criterion], cluster_assignments, mean)
496
497   if (should_pick_min(input$numeric_criterion)) {
498     recommended_cluster <- which.min(cluster_means)
499   } else {
500     recommended_cluster <- which.max(cluster_means)
501   }
502
503   recommendation <- filtered_data[cluster_assignments == recommended_cluster, ]
504
505   # Limit the recommendation to the top 45 rows
506   recommendation <- head(recommendation, 45)
507 } else {
508   clustering_model <- dbscan(combined_features, eps = input$dbscan_eps, minpts = input$dbscan_minpts)
509   cluster_means <- tapply(filtered_data[input$numeric_criterion][clustering_model$cluster != -1],
510     clustering_model$cluster[clustering_model$cluster != -1],
511     mean)
512
513   if (should_pick_min(input$numeric_criterion)) {
514     recommended_cluster <- which.min(cluster_means)
515   } else {
516     recommended_cluster <- which.max(cluster_means)
517   }
518
519   recommendation <- filtered_data[clustering_model$cluster == recommended_cluster, ]
520
521   # Limit the recommendation to the top 45 rows
522   recommendation <- head(recommendation, 45)
523 }
524
525 recommendation <- recommendation[order(recommendation$resale_price), ]
526 # Limit the recommendation to the top 45 rows
527 recommendation <- head(recommendation, 45)
528
529 if (nrow(recommendation) == 0) return("No recommendation available")
530

```

Figure C.5: unsupervised recommendation system shiny app part 5

```

533 ~ }
534 ~ output$filtered_data <- renderTable({
535 ~   recommendation_data()
536 ~ })
537 ~ if (nrow(recommendation) == 0) {
538 ~   return("No recommendation available")
539 ~ }
540 ~
541 ~ # Add a 'row_number' column to the start of your dataframe
542 ~ recommendation <- tibble::rowid_to_column(recommendation, "Row Number")
543 ~
544 ~ return(recommendation)
545 ~ })
546 ~
547 ~ output$map <- renderLeaflet({
548 ~   recommendation_data()
549 ~ })
550 ~ if (nrow(recommendation) == 0) {
551 ~   return(leaflet() %>% setView(lat = 1.3521, lng = 103.8198, zoom = 12)) # Default view of Singapore
552 ~ }
553 ~
554 ~ m <- leaflet() %>%
555 ~   addTiles() %>% # This adds the base map
556 ~   setView(lat = 1.3521, lng = 103.8198, zoom = 12) # Centered on Singapore
557 ~
558 ~ # Add a new column to 'recommendation' that represents the row number
559 ~ recommendation$row_number <- seq_len(nrow(recommendation))
560 ~
561 ~ # Use the 'popup' argument to display the row number along with any other info you want
562 ~ popups <- sprintf(
563 ~   "<strong>Row number:</strong> %s<br/><strong>Resale Price:</strong> %s",
564 ~   recommendation$row_number,
565 ~   recommendation$resale_price # Add other details as needed
566 ~ )
567 ~
568 ~ m <- addMarkers(
569 ~   map = m,
570 ~   data = recommendation,
571 ~   lat = ~Latitude,
572 ~   lng = ~Longitude,
573 ~   popup = popups # this argument accepts HTML content
574 ~ )
575 ~
576 ~ return(m)
577 ~ })
578 ~ }
579 ~
580 ~ shinyApp(ui, server)
581 ~

```

Figure C.6: unsupervised recommendation system shiny app part 6

```

566- '''[r training and testing model]
567- # Install and load necessary libraries
568- library(tidyverse)
569- library(lubridate)
570- library(caret)
571- library(randomForest)
572- library(xgboost)
573- library(e1071)
574- library(ggplot2)
575- library(gbm)
576- library(glmnet)
577- library(grid)
578-
579- # Load the dataset
580- data <- read_csv('final.csv', show_col_types = FALSE)
581-
582- # Convert categorical variables
583- data <- data %>%
584-   mutate_at(vars(town, flat_type, flat_model, storey_range), as.factor)
585-
586- # Convert 'remaining_lease' into numeric (in terms of years)
587- data$remaining_lease <- sapply(strsplit(as.character(data$remaining_lease), " "), function(x) {
588-   years <- as.numeric(x[1])
589-   if(length(x) > 2) {
590-     months <- as.numeric(gsub("month.", "", x[3])) / 12
591-     return(years + months)
592-   } else {
593-     return(years)
594-   }
595- })
596-
597- # Split data into training and test sets
598- set.seed(123)
599- trainIndex <- createDataPartition(data$resale_price, p = 0.8, list = FALSE)
600- train_data <- data[trainIndex,]
601- test_data <- data[-trainIndex,]
602-
603- # Train Random Forest
604- set.seed(123)
605- rf_model <- randomForest(resale_price ~ ., data=train_data, ntree=50)
606-
607- # Train XGBoost
608- xgb_data <- train_data %>%
609-   mutate_if(is.character, as.numeric) %>%
610-   mutate_if(is.factor, as.numeric)
611- xgb_train <- xgb.DMatrix(data = as.matrix(xgb_data %>% select(-resale_price)), label = xgb_data$resale_price)
612- set.seed(123)
613- xgb_model <- xgboost(data = xgb_train, nrounds=100, objective="reg:squarederror")
614-
615- # Evaluate models on test set
616- rf_predictions <- predict(rf_model, newdata = test_data)
617- xgb_predictions <- predict(xgb_model, newdata = as.matrix(test_data %>% select(-resale_price) %>% mutate_if(is.character, as.numeric) %>% mutate_if(is.factor, as.numeric)))
618-
619- rf_rmse <- sqrt(mean((rf_predictions - test_data$resale_price)^2))
620- xgb_rmse <- sqrt(mean((xgb_predictions - test_data$resale_price)^2))
621-
622- # SVM for regression
623- set.seed(123)
624- svm_model <- svm(resale_price ~ ., data=train_data, kernel="radial")
625- svm_predictions <- predict(svm_model, newdata = test_data)
626- svm_rmse <- sqrt(mean((svm_predictions - test_data$resale_price)^2))
627-
628- # 1. Gradient Boosting Machine (GBM)
629- set.seed(123)

```

Figure C.7: Training and testing model part 1

```

629 set.seed(123)
630 gbm_model <- gbm(resale_price ~ ., data=train_data, distribution="gaussian", n.trees=100)
631 gbm_predictions <- predict(gbm_model, newdata=test_data, n.trees=100)
632 gbm_rmse <- sqrt(mean((gbm_predictions - test_data$resale_price)^2))
633
634 # 2. LASSO Regression
635 # First, we need to create a model matrix
636 x_train <- model.matrix(resale_price ~ ., train_data)[,-1]
637 y_train <- train_data$resale_price
638 x_test <- model.matrix(resale_price ~ ., test_data)[,-1]
639 lasso_model <- cv.glmnet(x_train, y_train, alpha=1)
640 lasso_predictions <- predict(lasso_model, s=lasso_model$lambda.min, newx=x_test)
641 lasso_rmse <- sqrt(mean((lasso_predictions - test_data$resale_price)^2))
642
643 # 3. Ridge Regression
644 ridge_model <- cv.glmnet(x_train, y_train, alpha=0)
645 ridge_predictions <- predict(ridge_model, s=ridge_model$lambda.min, newx=x_test)
646 ridge_rmse <- sqrt(mean((ridge_predictions - test_data$resale_price)^2))
647
648 # Function to compute MAPE
649 calculate_mape <- function(actual, predicted) {
650   return(mean(abs((actual - predicted) / actual)) * 100)
651 }
652
653 # Calculate MAPE for each model
654 rf_mape <- calculate_mape(test_data$resale_price, rf_predictions)
655 xgb_mape <- calculate_mape(test_data$resale_price, xgb_predictions)
656 svm_mape <- calculate_mape(test_data$resale_price, svm_predictions)
657 gbm_mape <- calculate_mape(test_data$resale_price, gbm_predictions)
658 lasso_mape <- calculate_mape(test_data$resale_price, lasso_predictions)
659 ridge_mape <- calculate_mape(test_data$resale_price, ridge_predictions)
660
661 # Add MAPE to the comparison dataframe
662 comparison_df <- tibble(
663   Model = c("Random Forest", "XGBoost", "SVM", "GBM", "LASSO", "Ridge"),
664   RMSE = c(rf_rmse, xgb_rmse, svm_rmse, gbm_rmse, lasso_rmse, ridge_rmse),
665   MAPE = c(rf_mape, xgb_mape, svm_mape, gbm_mape, lasso_mape, ridge_mape)
666 )
667
668 print(comparison_df)
669 # Find the model with the lowest RMSE
670 best_model <- comparison_df[which.min(comparison_df$RMSE), ]
671
672 print(best_model)
673
674 # Plotting RMSE comparison
675 rmse_plot <- ggplot(comparison_df, aes(x = Model, y = RMSE)) +
676   geom_bar(stat = 'identity') +
677   coord_flip() +
678   labs(title = "RMSE Comparison across Models", x = "Models", y = "RMSE") +
679   theme_minimal() +
680   theme(legend.position = "none")
681
682 # Save the RMSE plot
683 ggsave(filename = "rmse_comparison_plot.png", plot = rmse_plot, width = 10, height = 6, bg = "white")
684
685 # Plotting MAPE comparison
686 mape_plot <- ggplot(comparison_df, aes(x = Model, y = MAPE)) +
687   geom_bar(stat = 'identity') +
688   coord_flip() +
689   labs(title = "MAPE Comparison across Models", x = "Models", y = "MAPE (%)") +
690   theme_minimal() +
691   theme(legend.position = "none")
692
693 # Save the MAPE plot
694 ggsave(filename = "mape_comparison_plot.png", plot = mape_plot, width = 10, height = 6, bg = "white")

```

Figure C.8: Training and testing model part 2

```

696 # Extract feature importance
697 rf_importance_matrix <- importance(rf_model)
698
699 # Create a data frame for feature importance using IncNodePurity
700 rf_importance <- data.frame(Feature = rownames(rf_importance_matrix), Importance = rf_importance_matrix[, "IncNodePurity"])
701
702 # Sort the importance data in descending order for better visualization
703 rf_importance <- rf_importance[order(-rf_importance$Importance), ]
704
705 # Save the plot with a white background
706 rf_plot <- ggplot(rf_importance, aes(x = reorder(Feature, Importance), y = Importance)) +
707   geom_bar(stat = 'identity') +
708   coord_flip() +
709   labs(title = "Feature Importance from Random Forest", x = "Features", y = "Importance (IncNodePurity)") +
710   theme_minimal() +
711   theme(axis.text.y = element_text(size = 12)) # Increase text size
712
713 # Save the plot with a white background
714 ggsave("rf_feature_importance.png", plot = rf_plot, width = 10, height = 12, bg = "white")
715
716 # Extract feature importance
717 xgb_importance <- xgb.importance(model = xgb_model)
718
719 # Saving the XGBoost feature importance plot
720 png(filename = "xgb_feature_importance.png", width = 800, height = 600)
721 xgb_plot <- ggplot(xgb_importance, aes(x = Feature, y = Importance)) +
722   geom_bar(stat = 'identity') +
723   labs(title = "Feature Importance (XGBoost)")
724 dev.off()
725
726 ...

```

Figure C.9: Training and testing model part 3

```

729 # manually predicting
730 predict_resale_price <- function(coun, flat_type, flat_model, storey_range,
731 floor_area_son, remaining_lease,
732 hamker_closest_distance, num_hamker_1km,
733 nrt_closest_dist, num_nrt_1km,
734 park_closest_dist, num_park_1km,
735 sch_closest_dist, num_sch_1km,
736 shoppingmall_closest_dist, num_shoppingmall_1km,
737 supermarket_closest_dist, num_supermarket_1km) {
738 # convert input parameters to a data frame
739 input_data <- as.data.frame(list(coun = coun, flat_type = flat_type,
740 flat_model = flat_model, storey_range = storey_range,
741 floor_area_son = floor_area_son, remaining_lease = remaining_lease,
742 hamker_closest_distance = hamker_closest_distance,
743 num_hamker_1km = num_hamker_1km,
744 nrt_closest_dist = nrt_closest_dist, num_nrt_1km = num_nrt_1km,
745 park_closest_dist = park_closest_dist, num_park_1km = num_park_1km,
746 sch_closest_dist = sch_closest_dist, num_sch_1km = num_sch_1km,
747 shoppingmall_closest_dist = shoppingmall_closest_dist,
748 num_shoppingmall_1km = num_shoppingmall_1km,
749 supermarket_closest_dist = supermarket_closest_dist,
750 num_supermarket_1km = num_supermarket_1km))
751 # convert categorical variables to factor
752 input_data[["coun"]] <- factor(input_data[["coun"]], as.factor())
753 input_data[["flat_type"]] <- factor(input_data[["flat_type"]], as.factor())
754 input_data[["flat_model"]] <- factor(input_data[["flat_model"]], as.factor())
755 input_data[["storey_range"]] <- factor(input_data[["storey_range"]], as.factor())
756 # Adjust the levels of the factors in input_data to match those in the original data
757 input_data[["coun"]] <- factor(input_data[["coun"]], levels = levels(data[["coun"]]))
758 input_data[["flat_type"]] <- factor(input_data[["flat_type"]], levels = levels(data[["flat_type"]]))
759 input_data[["flat_model"]] <- factor(input_data[["flat_model"]], levels = levels(data[["flat_model"]]))
760 input_data[["storey_range"]] <- factor(input_data[["storey_range"]], levels = levels(data[["storey_range"]]))
761 # convert the newdata for xgboost prediction
762 xgb_input_data <- input_data[["coun"]]
763 xgb_input_data[["flat_type"]] <- input_data[["flat_type"]]
764 xgb_input_data[["flat_model"]] <- input_data[["flat_model"]]
765 xgb_input_data[["storey_range"]] <- input_data[["storey_range"]]
766 # ensure that the column order matches the training data for xgboost
767 xgb_input_data <- xgb_input_data[, colnames(xgb_data) %in% select(resale_price)]
768 # make predictions
769 rf_pred <- unname(predict(rf_model, newdata = input_data)) # unname the prediction
770 xgb_pred <- predict(xgb_model, newdata = as.matrix(xgb_input_data))
771 # calculate range
772 rf_range <- unname(c(rf_pred * (1 - range_percentage), rf_pred * (1 + range_percentage))) # unname the range
773 xgb_range <- c(xgb_pred * (1 - range_percentage), xgb_pred * (1 + range_percentage))
774 return(list(
775 rf_prediction = rf_pred,
776 rf_range = rf_range,
777 xgb_prediction = xgb_pred,
778 xgb_range = xgb_range
779 ))
780 }
781
782 result <- predict_resale_price(
783 coun = "sg",
784 flat_type = "2 room",
785 flat_model = "improved",
786 storey_range = "07 to 09",
787 floor_area_son = 44, # 12.22, # 38 years and 1 month
788 remaining_lease = 18, # 12.22, # 38 years and 1 month
789 hamker_closest_distance = 0.100276779,
790 num_hamker_1km = 0,
791 nrt_closest_dist = 0.999999999,
792 num_nrt_1km = 0,
793 park_closest_dist = 0.674068937,
794 num_park_1km = 21,
795 sch_closest_dist = 0.333306119,
796 num_sch_1km = 30,
797 shoppingmall_closest_dist = 1.002781394,
798 num_shoppingmall_1km = 6,
799 supermarket_closest_dist = 0.742126182,
800 num_supermarket_1km = 18
801 )
802 print(result)
803
804 # ...

```

Figure C.10: Manually predicting housing resale price based on the parameters

```

828 # ... for shiny app
829 library(shiny)
830 library(shinydashboard)
831 library(shinyxgboost)
832 # define ui for application
833 ui <- fluidPage(
834 # application title
835 titlePanel("Predict Resale Price"),
836 # sidebar with input elements
837 sidebarLayout(
838 # sidebar
839 selectInput("coun", "Coun", choices = levels(data[["coun"]])),
840 selectInput("flat_type", "Flat Type", choices = levels(data[["flat_type"]])),
841 selectInput("flat_model", "Flat Model", choices = levels(data[["flat_model"]])),
842 selectInput("storey_range", "Storey Range", choices = levels(data[["storey_range"]])),
843 numericInput("floor_area_son", "Floor Area (sqm)", 44),
844 numericInput("remaining_lease", "Remaining Lease (in years)", 18),
845 numericInput("hamker_closest_distance", "Closest Hamker Distance", 0.10),
846 numericInput("num_hamker_1km", "Number of Hamkers within 1km", 0),
847 numericInput("nrt_closest_dist", "Closest NRT Distance", 0.9999),
848 numericInput("num_nrt_1km", "Number of NRTs within 1km", 0),
849 numericInput("park_closest_dist", "Closest Park Distance", 0.674),
850 numericInput("num_park_1km", "Number of Parks within 1km", 21),
851 numericInput("sch_closest_dist", "Closest School Distance", 0.333),
852 numericInput("num_sch_1km", "Number of Schools within 1km", 30),
853 numericInput("shoppingmall_closest_dist", "Closest Shopping Mall Distance", 1.003),
854 numericInput("num_shoppingmall_1km", "Number of Shopping Malls within 1km", 6),
855 numericInput("supermarket_closest_dist", "Closest Supermarket Distance", 0.742),
856 numericInput("num_supermarket_1km", "Number of Supermarkets within 1km", 18),
857 actionButton("predict", "Predict")
858 ),
859 # main panel for displaying predictions
860 mainPanel(
861 # random forest prediction
862 verbatimTextOutput("rf_prediction"),
863 # xgboost prediction
864 verbatimTextOutput("xgb_prediction")
865 )
866 )
867 # define server logic
868 server <- function(input, output) {
869 # observe variables for prediction
870 result <- predict_resale_price(
871 coun = input$coun,
872 flat_type = input$flat_type,
873 flat_model = input$flat_model,
874 storey_range = input$storey_range,
875 floor_area_son = input$floor_area_son,
876 remaining_lease = input$remaining_lease,
877 hamker_closest_distance = input$hamker_closest_distance,
878 num_hamker_1km = input$num_hamker_1km,
879 nrt_closest_dist = input$nrt_closest_dist,
880 num_nrt_1km = input$num_nrt_1km,
881 park_closest_dist = input$park_closest_dist,
882 num_park_1km = input$num_park_1km,
883 sch_closest_dist = input$sch_closest_dist,
884 num_sch_1km = input$num_sch_1km,
885 shoppingmall_closest_dist = input$shoppingmall_closest_dist,
886 num_shoppingmall_1km = input$num_shoppingmall_1km,
887 supermarket_closest_dist = input$supermarket_closest_dist,
888 num_supermarket_1km = input$num_supermarket_1km
889 )
890 output$rf_prediction <- renderPrint({
891 textOutput("rf_prediction", range = result$rf_range)
892 })
893 output$xgb_prediction <- renderPrint({
894 textOutput("xgb_prediction", range = result$xgb_range)
895 })
896 }
897 # run the application
898 shinyApp(ui = ui, server = server)
899
900 # ...

```

Figure C.11: Supervised predicting housing resale price shiny app