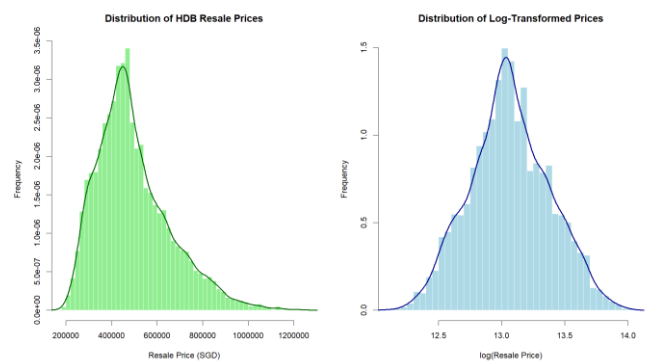


## **Part 0: Introduction**

The objective of this report is to develop and propose the best possible linear regression model to predict the resale price of Housing and Development Board (HDB) flats in Singapore. Understanding the factors that influence resale prices is important for assessing housing affordability and property market trends in a nation where over 80% of residents live in HDB flats. The dataset, obtained from [Data.gov.sg](https://data.gov.sg), contains 11 variables and 11,527 observations representing resale transactions between January and June 2021. Each record includes details such as the month of sale, town, flat type, floor area, flat model, storey range, and lease commencement year, alongside the resale price in Singapore dollars. These variables capture both physical attributes (e.g., flat size, storey height, architectural model) and locational characteristics (e.g., town or region). The analysis proceeds in several stages. Part I (EDA) explores the response and explanatory variables using summary statistics and visualizations such as histograms, scatter plots, and box plots to understand data patterns and relationships. Part II fits and refines multiple linear regression models, starting from an initial model ( $M_0$ ) with at least six regressors, to test the adequacy of the model, which includes the linearity, normality, and constant-variance assumptions, and to address issues such as multicollinearity and overfitting. Subsequent sections employ model selection criteria (Adjusted  $R^2$ , AIC, and plots) to identify the most adequate model. The report concludes with the presentation and interpretation of the final model ( $M_n$ ) and its practical implications for understanding the key determinants of HDB resale prices in Singapore.

## **Part I EDA: Exploring the variables and association**



*Figure 1.1*

The histogram in Figure 1.1 shows that HDB resale prices from January to June 2021 are right-skewed, with most transactions between SGD 300,000 and 1,000,000, and a few high-value flats extending the tail beyond 1 million. The overlaid density curve indicates clear deviation from normality, suggesting that using raw prices may violate linear regression assumptions. To correct this, the response variable was log-transformed, yielding a roughly symmetric, bell-shaped distribution (Figure 1.1). This transformation effectively reduces skewness and stabilizes variance, making  $\log(\text{resale price})$  a more appropriate response variable for subsequent regression modeling.

To prepare the data for analysis, storey range was converted from a categorical predictor (e.g., “07 TO 09”) to numeric midpoints (e.g., 8), allowing storey height to be treated as a quantitative predictor. The 26 towns were also grouped into three regional categories, Core Central Region (CCR), Rest of Central Region (RCR), and Outside Central Region (OCR), to reduce dimensionality and enhance interpretability of location effects, as well as to prevent excessive coefficients in our models. Variables such as block and street name were removed because they capture overly specific location details that do not generalize well. Since the regional grouping already represents geographic variation, including these variables would add redundancy and risk multicollinearity without improving explanatory power.

In summary, the dataset comprises both quantitative variables (e.g., floor area, storey range, lease commencement date, remaining lease) and categorical variables (e.g., month, flat type, flat model, region), allowing for a comprehensive analysis of how both continuous and qualitative attributes influence HDB resale prices

The scatterplots in Figure 1.2 illustrate the relationships between  $\log(\text{resale price})$  and four quantitative predictors: floor area, lease commencement date, storey range, and remaining lease. Floor area shows the strongest positive linear relationship, indicating that larger flats generally command higher resale prices. A moderate positive trend is also observed for storey height, as higher floors tend to be valued more for better views and ventilation. Lease commencement date and remaining lease both exhibit mild positive relationships with price, reflecting the depreciation effect of Singapore's 99-year leasehold system, as newer flats or those with longer leases sell for more. Although some dispersion and mild heteroscedasticity remain, the overall trends are approximately linear, making these predictors suitable for multiple linear regression analysis.

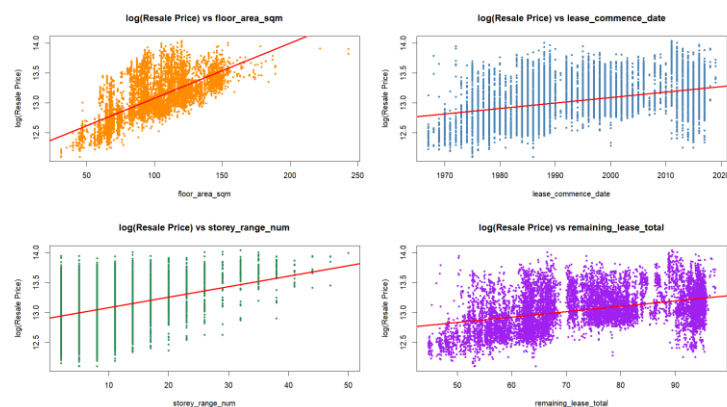


Figure 1.2

The boxplots in Figure 1.3 below show variations in  $\log(\text{resale price})$  across the categorical variables: month, flat type, flat model, and region. Prices remain relatively stable across months from January to June 2021, indicating minimal temporal effects within this period. In contrast, flat type displays a clear upward trend, where larger units (4-room, 5-room, Executive, and Multi-Generation) command higher resale prices, consistent with the earlier positive association between floor area and price. The flat model variable also shows notable price variation across architectural designs, older models such as Improved or Standard have lower median prices, while premium designs like Maisonette and Terrace achieve higher valuations, capturing structural differences not explained by size alone. Finally, region exhibits the expected location gradient: flats in the Core Central Region (CCR) are most expensive, followed by the Rest of Central Region (RCR) and Outside Central Region (OCR). Overall, flat type, flat model, and region emerge as key categorical predictors of resale price, while month plays a minor role during this six-month window.

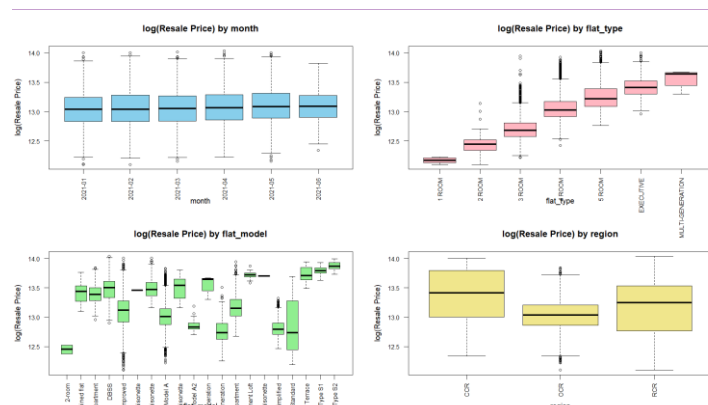


Figure 1.3

Some predictors in the dataset are likely interrelated. For example, floor area correlates with flat type, as larger flats typically have more rooms, while remaining lease is closely tied to lease commencement date, since newer flats naturally retain longer leases. These relationships suggest potential multicollinearity, where overlapping information among predictors may distort coefficient estimates. This issue will be formally assessed using Variance Inflation Factors (VIF) to identify and mitigate redundancy, ensuring a more stable and parsimonious model. Additionally, model selection techniques such as AIC, BIC, and stepwise regression will be employed to determine an optimal subset of variables that balances goodness of fit with model simplicity, ensuring that the final model is both parsimonious and statistically robust.

## Part II Building Model:

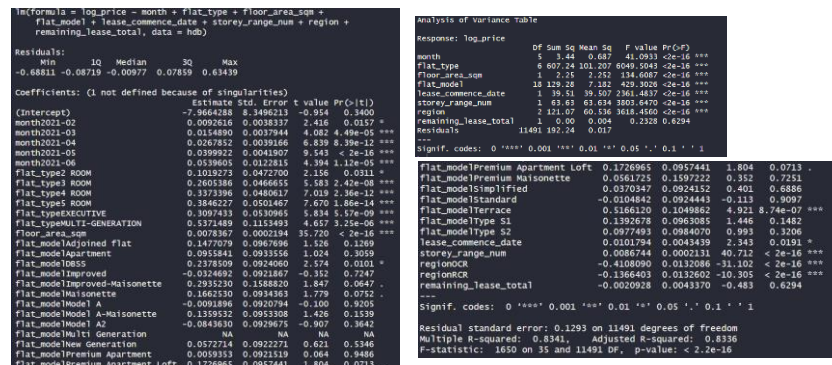


Figure 2.1

**Initial Model (M<sub>0</sub>):** Figure 2.1 shows that the initial multiple linear regression model (M<sub>0</sub>) was fitted using all major categorical and quantitative predictors. M<sub>0</sub> explains a high proportion of the variance in log(resale price) ( $R^2 = 0.8341$ , Adjusted  $R^2 = 0.8336$ ), with the overall F-statistic ( $p < 2.2e-16$ ) confirming strong model significance. Among predictors, floor area and storey range are highly significant ( $p < 0.001$ ), indicating that larger flats and higher floors are associated with higher resale prices, while lease commencement date also has a positive effect ( $p = 0.0191$ ). In contrast, several levels of flat model show large p-values, suggesting redundancy or possible overlap with flat type and floor area. The ANOVA results reinforce that flat type, region, and storey range explain substantial variance, whereas remaining lease total contributes little, as its effect is already captured by lease commencement date. Overall, M<sub>0</sub> demonstrates strong explanatory power but exhibits signs of over-parameterization and multicollinearity, motivating refinement through model simplification and diagnostic analysis.

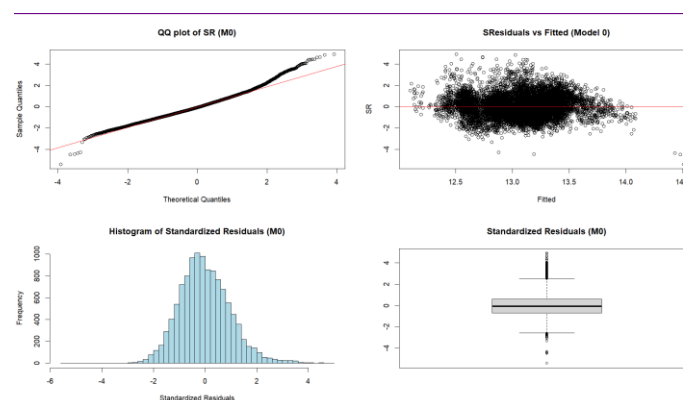


Figure 2.2

**Model Adequacy Check:** Model adequacy for M<sub>0</sub> was evaluated using diagnostic plots of standardized residuals as shown in Figure 2.2. The QQ plot and histogram indicate approximate

normality, with only slight tail deviations likely caused by a few outliers ( $\approx 203$  cases out of 11,527) as shown in the boxplot, which have minimal impact on model fit. The residuals-versus-fitted y-values plot also show generally linear trend, suggesting the linearity assumption is reasonably satisfied. However, the residuals-versus-fitted y-values plot displays a clear funnel shape, indicating violation of the constant variance assumption. Overall, while linearity and normality hold adequately, the non-constant variance suggests the need for transformation or a Weighted Least Squares (WLS) approach to improve model adequacy.

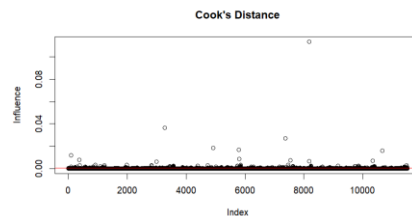


Figure 2.3

**Check for Influential Points:** In Figure 2.3 shown, the Cook's Distance plot was used to detect influential observations that might disproportionately affect the model. Almost all points have values near zero, and none exceed the conventional threshold of 1, indicating no influential observations. Although a few data points show slightly higher values, they do not significantly impact the regression results. Thus, while  $M_0$  contains minor outliers, there is no evidence of undue influence compromising model stability or reliability.

```
> vif_vals
```

	VIF	Df	GVIFA(1/(2*Df))
month	1.217155	5	1.019846
flat_type	268.960155	5	1.749721
Floor_area_sqm	18.654608	1	4.319121
flat_model	115.838962	18	1.141115
lease_commence_date	2691.134285	1	51.876144
storey_range_num	1.238637	1	1.112941
region	1.878900	2	1.170782
remaining_lease_total	2695.136244	1	51.914702

Figure 2.4

**Check for Multicollinearity using VIF:** As shown in Figure 2.4, to evaluate multicollinearity, the Variance Inflation Factor (VIF) was computed for all predictors. Most variables showed low VIF values ( $< 5$ ), indicating minimal collinearity. However, lease commence date and remaining lease total displayed extremely high adjusted GVIFs ( $\approx 52$ ), reflecting severe redundancy since both capture similar lease-age information. Consequently, remaining lease total was removed through backward stepwise regression based on the Akaike Information Criterion (AIC), resulting in a refined model ( $M_1$ ) with a slightly lower AIC and marginally improved adjusted  $R^2$ . This confirms that eliminating redundant variables enhances parsimony without sacrificing explanatory power.

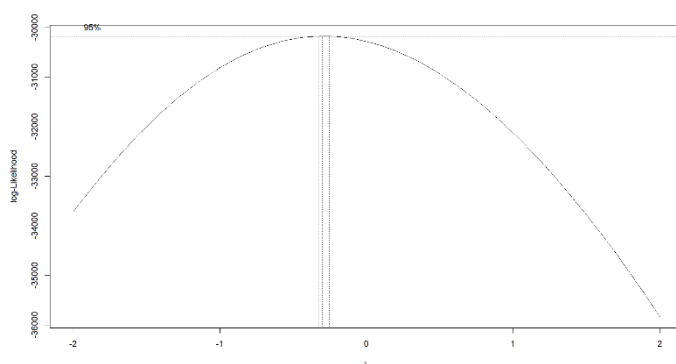


Figure 2.5

```
> AIC(M2, M3, M4)
```

	df	AIC
M2	35	-14528.99
M3	35	-14514.61
M4	36	-14280.84

Figure 2.6

**Transformations:** To confirm the suitability of transforming the response variable, a Box–Cox transformation was applied with resale price as the dependent variable. The resulting plot showed an optimal  $\lambda$  near zero, with the 95% confidence interval tightly surrounding  $\lambda = 0$ , validating the use of  $\log(\text{resale price})$  as the response (Figure 2.5). This transformation effectively stabilized variance and improved residual normality. To further address mild non-linearity and heteroscedasticity, floor area per square metre and lease commence date were also log-transformed, producing Model M<sub>2</sub>, which achieved an improved fit (AIC = −14528.99) (Figure 2.6). A Weighted Least Squares (WLS) model (M<sub>3</sub>) was then fitted using weights inversely proportional to the squared fitted values from M<sub>2</sub>, reducing variance non-constancy though yielding a slightly higher AIC (−14514.61) (Figure 2.6). A subsequent Generalized WLS model (M<sub>4</sub>) with a varPower variance structure produced a much higher AIC (−14280.84) (Figure 2.6), indicating that its added complexity did not enhance model performance. Hence, M<sub>4</sub> was rejected due to over-parameterization and inferior fit.

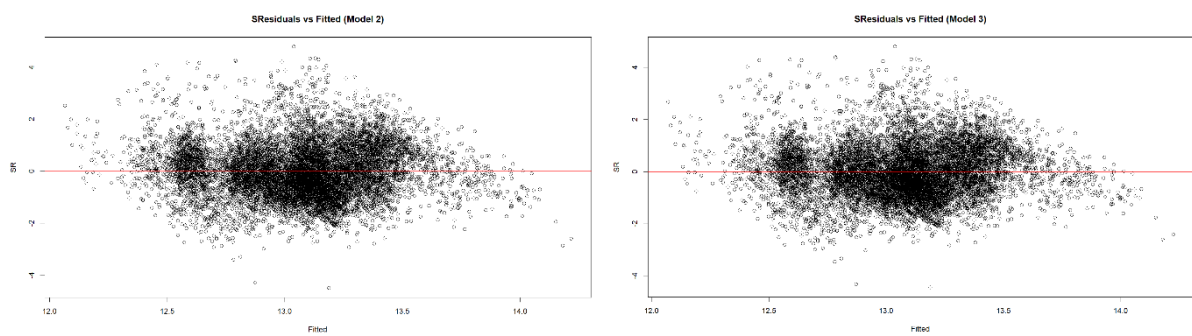


Figure 2.7

**Check for Correction of Inadequacy:** As shown in Figure 2.7, residual plots for Model 2 (OLS) and Model 3 (WLS) were compared to assess the constant-variance assumption. In M<sub>2</sub>, standardized residuals show a mild funnel pattern, indicating slight heteroscedasticity despite the earlier log transformations. Applying inverse-squared fitted value weights in M<sub>3</sub> produced a more uniform residual spread, substantially reducing this pattern and satisfying the homoscedasticity assumption. Although M<sub>3</sub> has a slightly higher AIC (−14514.61 vs −14528.99), its improved variance stability makes it statistically more adequate overall.

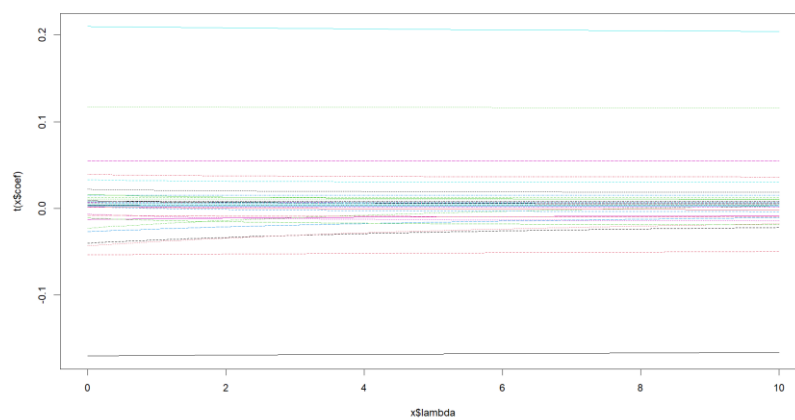


Figure 2.8

**Check for Multicollinearity using Ridge Regression:** To evaluate multicollinearity in the final weighted model (M<sub>3</sub>), ridge regression diagnostics were performed. The ridge trace plot shows



stable, nearly flat coefficient paths across penalty values ( $\lambda = 0-10$ ), with no major fluctuations or crossovers (Figure 2.8). This indicates that the coefficient estimates are stable and multicollinearity is minimal within  $M_3$ . Hence, since  $M_3$  already achieves low collinearity, regularization was unnecessary.

### Chosen Model ( $M_3$ ):

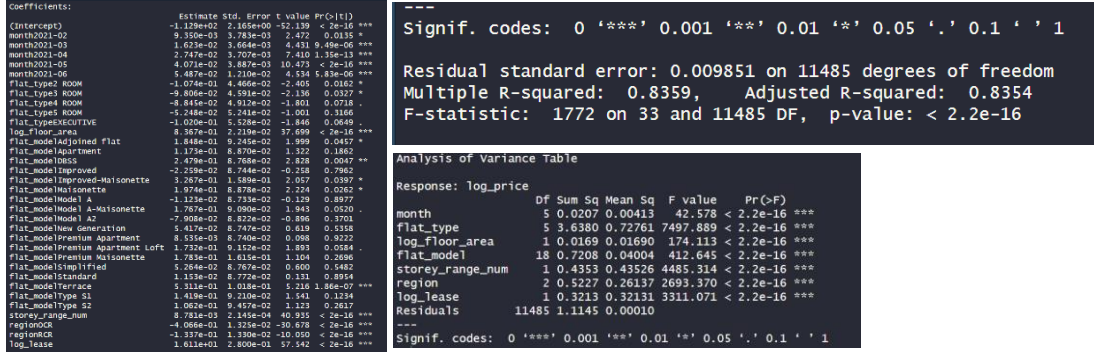


Figure 3

**Conclusion:** Model 3 (WLS) demonstrates the most adequate performance among all competing models. The application of weighted least squares effectively stabilized residual variance, as shown by the residual–fitted plot, which no longer exhibits the funnel pattern seen in  $M_2$ . Although its AIC (–14514.61) is slightly higher than  $M_2$ 's (–14528.99) (Figure 2.6), this small trade-off is justified by superior compliance with the constant-variance assumption. The model also avoids over-parameterization, maintaining a balanced set of predictors that each contribute meaningfully to explaining variation in log(resale price). Statistically,  $M_3$  exhibits strong explanatory power ( $R^2 = 0.8359$ , Adjusted  $R^2 = 0.8354$ ), with the overall F-statistic ( $F = 1772$ ,  $p < 2.2e-16$ ) confirming high model significance (Figure 3). The residual standard error of 0.00985 further indicates excellent model precision. Overall,  $M_3$  achieves both strong explanatory performance and full compliance with regression assumptions, linearity, normality, homoscedasticity, and low multicollinearity, making it the most statistically robust and practically interpretable model for predicting HDB resale prices.

In practical terms, the final model highlights how physical characteristics (flat size, storey height, design type) and locational attributes (region and lease recency) jointly shape HDB resale values in Singapore. These findings reinforce that newer, larger, and more centrally located flats command significant price premiums. However, the model is limited to transactions from the first half of 2021 and does not account for external factors such as proximity to amenities, renovations, or macroeconomic changes. Future analyses could incorporate these variables or extend the dataset across multiple years to enhance predictive generalizability.

### Final Model Equation:

$$\widehat{\log(\text{Resale Price})} = \hat{\beta}_0 + \sum_{i=1}^5 \hat{\beta}_{1i}(\text{Month}_i) + \sum_{j=1}^6 \hat{\beta}_{2j}(\text{Flat Type}_j) + \hat{\beta}_3 \log(\text{Floor Area})$$

$$+ \sum_{k=1}^{18} \hat{\beta}_{4k}(\text{Flat Model}_k) + \hat{\beta}_5(\text{Storey Range}_{\text{num}}) + \sum_{m=1}^2 \hat{\beta}_{6m}(\text{Region}_m) + \hat{\beta}_7 \log(\text{Lease})$$