

The dataset that I have selected is HDB resale prices 2021. This data is taken from open sources (Data.gov.sg). This dataset is of personal interest to me as my mom works at HDB so this was an excellent opportunity for me to play with data related to her workplace. Additionally, I find it fascinating how I can use the concepts that I learned in DSE1101 and apply them to real-world data. Predicting resale prices involves building predictive models, which is a fundamental skill in data science. The models that I will be using are Decision Trees, Multiple Linear Regression, and K-nearest Neighbours (KNN). Since this is a regression problem, the models are evaluated based on the out-of-sample mean squared error (MSE) using a test set.

Data Preparation:

I checked for missing values and handled them by omitting the respective rows. The resale prices are scaled by dividing them by 1000 for ease of interpretation. I performed a 50/50 train-test split, with 3000 observations in the training set and the rest in the testing set. I also set a seed for reproducibility.

Variable Selection:

I first deduced a few variables that would be important in my study. The variable that I felt was the most important is the “floor_area_sqm” variable. This is a no-brainer as floor area is the first thing that shows up in every advertising brochure/website pertaining to a resale flat. It is so crucial that many resale websites such as propertyguru.com.sg have indicated a dollar-per-floor area metric, for buyers to judge the value of the house. Based on my understanding, with a larger floor area, homeowners would have a larger area to play around with, in terms of interior design and furniture and therefore, house valuation would be higher. Based on my economics knowledge, I also considered “Remaining_lease” and “Dist_nearest_station” to be crucial as well. When deciding which HDB to buy, rational buyers would consider the opportunity cost in making the decision. They would prefer to choose flats closer to MRT stations to cut down the implicit cost incurred which is the time taken to walk to the nearest MRT station. Similarly, they would choose HDBs with a higher remaining lease value so that they can cut down the time taken to search for a new flat.

Additionally, I used Decision Trees for variable selection. It is a non-parametric method which can perform regression tasks. I choose trees for many reasons. It is computationally efficient, has automatic interaction detection, handles both categorical and numeric X well and performs variable selection. I first fitted the entire data set and generated a big tree, using a small value $\text{mindev}=0.0001$ as a stopping criterion. The optimal tree size (13 leaves) was determined through 10-fold cross-validation. The original tree was pruned to the optimal size and the pruned tree is visualised (as shown in Figure 1) and assessed on the test data. Ultimately, the out-of-sample MSE is 6496.478.

```

graph TD
    Root[floor_area_sqm < 82.5] -->|Yes| Node1[max_floor_lvl < 24.5]
    Root -->|No| Node2[Dist_CBD < 8.11344]
    Node1 -->|Yes| Leaf1[324.8]
    Node1 -->|No| Leaf2[553.9]
    Node2 -->|Yes| Node3[max_floor_lvl < 27.5]
    Node2 -->|No| Node4[floor_area_sqm < 111.5]
    Node3 -->|Yes| Leaf3[836.4]
    Node3 -->|No| Node5[floor_area_sqm < 109]
    Node4 -->|Yes| Node6[flat_model_dbss < 0.5]
    Node4 -->|No| Node7[floor_area_sqm < 137.5]
    Node5 -->|Yes| Node8[Remaining_lease < 73.5]
    Node5 -->|No| Leaf4[742.9]
    Node6 -->|Yes| Node9[Remaining_lease < 82.5]
    Node6 -->|No| Leaf5[693.4]
    Node7 -->|Yes| Node10[Remaining_lease < 87.5]
    Node7 -->|No| Leaf6[694.0]
    Node8 -->|Yes| Leaf7[497.2]
    Node8 -->|No| Leaf8[676.2]
    Node9 -->|Yes| Leaf9[422.5]
    Node9 -->|No| Node11[Dist_CBD < 10.7188]
    Node10 -->|Yes| Leaf10[525.7]
    Node10 -->|No| Leaf11[625.4]
    Node11 -->|Yes| Leaf12[666.6]
    Node11 -->|No| Leaf13[480.6]
  
```

1. floor_area_sqm 2. max_floor_lvl 3. Dist_CBD 4. Remaining_lease 5. flat_model_dbss

The HDB_resale_readme.txt file in the HDB_resale_prices folder mentioned that the main variables are:

Model fitting:

As mentioned previously, it was fitted using all the variables in the train data.

Multiple Linear regression:

The red line in the diagnostic plots was generally horizontal, indicating linearity.

(2) Fitted using the variables that I initially thought were crucial based on domain knowledge using the train data: (floor area sqm, Remaining lease, Dist nearest station)

The out-of-sample MSE is **13697.36**.

(3) As an exploration exercise, I went to find out what factors affected the market value of BTO flats in Singapore. From this HDB website* I gathered a few factors: Proximity to city or town centre (Dist_CBD), Accessibility to key transport nodes(Dist_nearest_station), Accessibility to key amenities, such as supermarkets, food centres or malls(Dist_nearest_mall), Larger flat sizes(floor_area_sqm) and Higher storey heights(max_floor_lvl). Logically these factors should also affect the resale price of HDBs and I wanted to test that. I mapped those factors to their corresponding variables in my dataset and I fitted those variables using the train data.

The out-of-sample MSE is **6885.307**.

(4) Lastly, I fitted all the variables in the train data to compare their performance to my model.

The out-of-sample MSE is **425400.6**.

K-nearest Neighbours (KNN)

(1) Fitted using the selected variables from the tree model on train data. The optimal $K = 4$ is determined through cross-validation (as shown in Figure 2b) The final KNN model is evaluated on the test set.

The out-of-sample MSE is **3129.765**.

(2) Lastly, I fitted all the variables in the train data to compare their performance to my model.

The out-of-sample MSE is **5818.863**.

Figure 2a

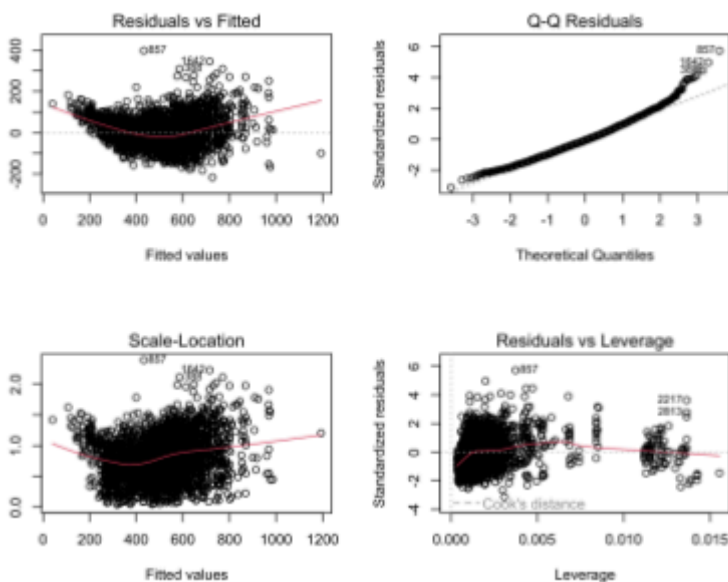
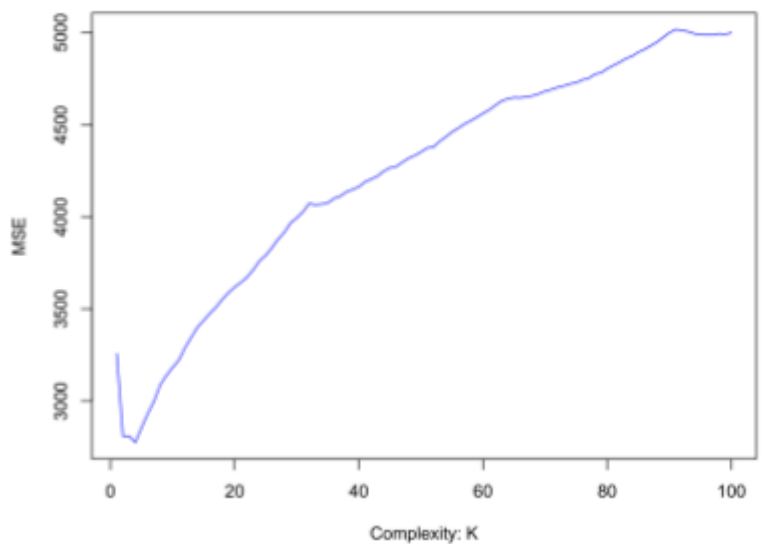


Figure 2b
LOOCV MSE



Analysis and limitations:

Overall, KNN with a K size of 4 had the lowest out-of-sample MSE of **3129.765**, followed by Logit with **4964.732**. This indicates KNN's superior predictive performance over Logit and Decision Trees. In Week

8's Lecture, KNN performed slightly worse than logit, but in my case, KNN performed better. However, that was a classification, not a regression problem and I did not use 1 fewer predictor this time. Both KNN and Logit performed better using the variables that I have selected using trees compared to fitting all the variables. This is because, upon analysis of the Logit(1) summary, all the t-values of the variables are more than 1.96 or less than -1.96, indicating they are all statistically significant and thus helpful for my analysis.

Limitations include:

1. Trees usually lose out to other methods of prediction as they are highly dependent on training data. A small change in data can result in a large change in the estimated tree. This would ultimately affect the variables that I chose for my model.
2. For Logit(3): I ignored factors such as "More ideal orientation with desirable views" and "Prevailing market conditions" from the HDB website as there were no variables in my dataset that I could map those factors to. These factors are not included in my analysis and could improve MSE, which was otherwise not represented.

As an exercise, I used my Logit(1) model to predict the price of a resale flat: 418 Serangoon Central**. The inputs are as follows: floor_area_sqm = 146.0436, max_floor_lvl = 12(I live nearby so I checked that the max floor is 12), Dist_CBD = 9.4(Used the shortest distance to Downtown Core on Google Maps), Remaining_lease = 64(Used HDB Emap*** to input the address and find remaining lease), flat_model_dbss = 0(Not a DBSS flat). The predicted price is \$705,932.2, whereas the actual price is \$1,150,000. This discrepancy could be due to benefits not accounted for such as the fact that it is a maisonette and its closeness to schools, transport and amenities.

Conclusion

It was interesting to find out how in Logit(2), the initial variables that I chose turned out to not perform well in fitting the model. This made me reflect on my understanding and made me read up on the many other factors that affect resale prices. These insights gained will aid in my decision-making when buying a resale flat in the future. Overall, this was an eye-opening experience to be able to explore handling large datasets using the models that were taught in class.

References:

*<https://www.hdb.gov.sg/about-us/news-and-publications/publications/hdbspeaks/How-BTO-Flats-are-Priced#:~:text=First%2C%20HDB%20establishes%20the%20market,individual%20attributes%20of%20the%20flats.&text=Projects%20with%20locational%20advantages%20will,to%20city%20or%20town%20centre>

**<https://www.propertyguru.com.sg/listing/hdb-for-sale-418-serangoon-central-24087219>

***<https://services2.hdb.gov.sg/web/fi10/emap.html>