

# **PROJECT 2: YOUR HOUSE, YOUR FUTURE: MAKING INFORMED REAL ESTATE DECISIONS**

**TEAM: DATA-NINE-NINE**

**Ming Fatt, Jasmine, Jin Jin, Wen Xi, Willson**

## DID YOU KNOW?



Expected Prices Increase of HDB

**2%- 8%** for year 2023



Whereas Price Hike of HDB

**10.4%** for year 2021

**12.7%** for year 2022



Source: <https://www.channelnewsasia.com/singapore/cooling-measures-singapore-hdb-resale-prices-towns-property-map-3499961#:~:text=Analysts%20expect%20a%20one%2Ddigit,12.7%20per%20cent%20in%202021>

# Problem Statement

The general public may not be well-equipped with the information needed to aid their **Real Estate** decision making process. Some of the common questions they might have are

- (1) What are the available options given my current budget?
- (2) Which flat type and where can I afford?
- (3) What price should I set when I sell my flat?
- (4) How to market my flat to increase its selling price?



# FRET NOT!

**REAL ESTATE START-UP COMPANY,  
DATA NINE-NINE IS HERE TO HELP!**

With our state-of-the-art data driven HDB resale price prediction model, your real estate issues shall be a thing of the past.





# **FLOW OF MODEL BUILDING PROCESS**

Understanding the model building process

# OUR PROCESS IS EASY



## Define the problem

- Identify market need
- Serve the need
- Through **Automated valuation process**

## Gather data

- Gather the necessary raw data
- Data cleaning works
- Feature Engineering

## Explore data

- Study correlation between features
- Verify reliability of correlation
- **Select features** for our model

# OUR PROCESS IS EASY



## Produce a model with the data

- Construct Model
- **Linear, Ridge, and Lasso Regression models**
- Optimise models

## Evaluate the model

- Evaluate models with Cross Validation, RMSE and R Squared
- Best model among the **9** are deployed

## Providing recommendation to the problem with the model

- Made **customised Predictions**
- Data-Driven Recommendations will be provided



# EXPLORATORY DATA ANALYSIS PROCESS

Data Janitorial Work

Datasets used contained

**150,634**

HDB flat resale transaction

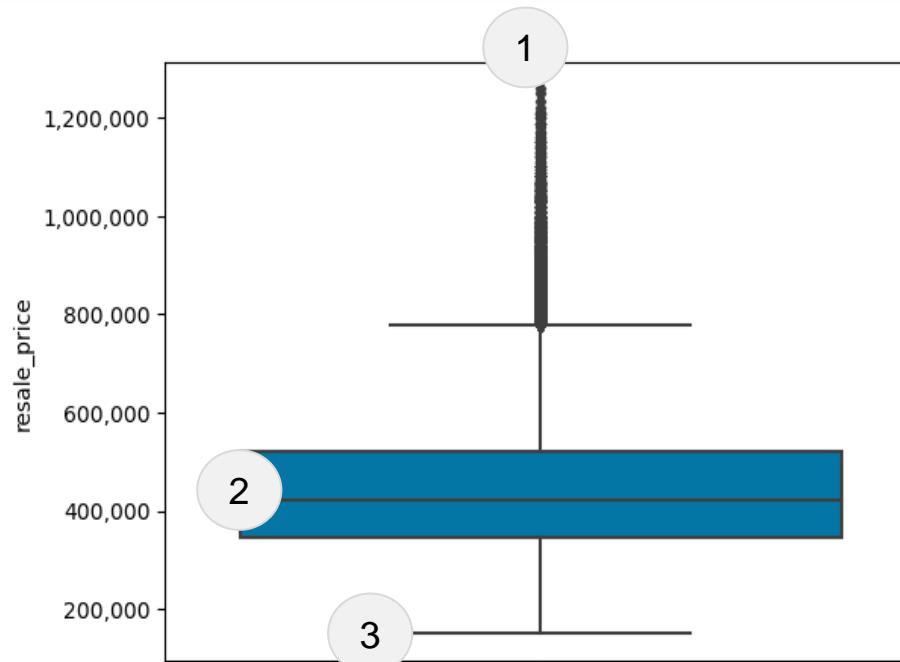
**77**

Features

**Mar '12 - Apr '21**

Duration of dataset

# Overview of sale pricing



1

**Most expensive sale**

\$1,258,000

5 room flat in Central Area  
(in 2020)

2

**Avg sale price**

\$449,162

3

**Cheapest sale**

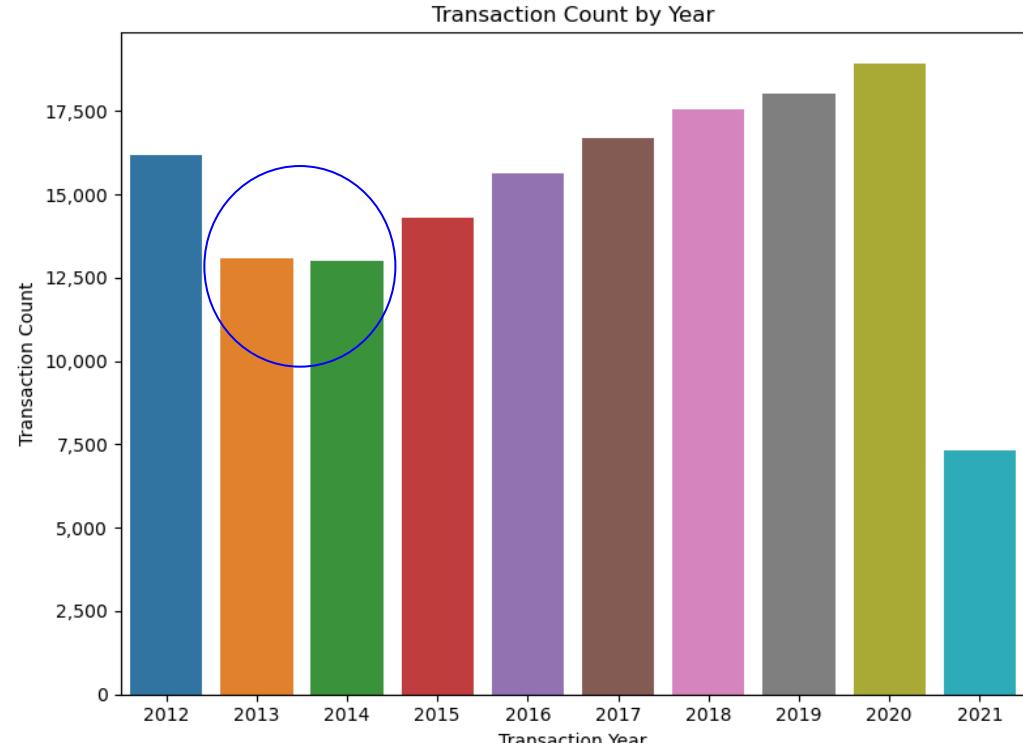
\$150,000

2 room flat in Toa Payoh  
(2020) & Geylang (in 2019)

# Exploratory Data Analysis



- The dip of transaction count of HDB resale flats in the year of 2013 & 2014 is noted, and may be related to cooling measures introduced in 2013\*.
- 2021 data is only up till Apr; if pro-rated for 2021(entire year) it is on track to be higher than in 2020



\* <https://stackedhomes.com/editorial/singapore-cooling-measures-history>

# Exploratory Data Analysis



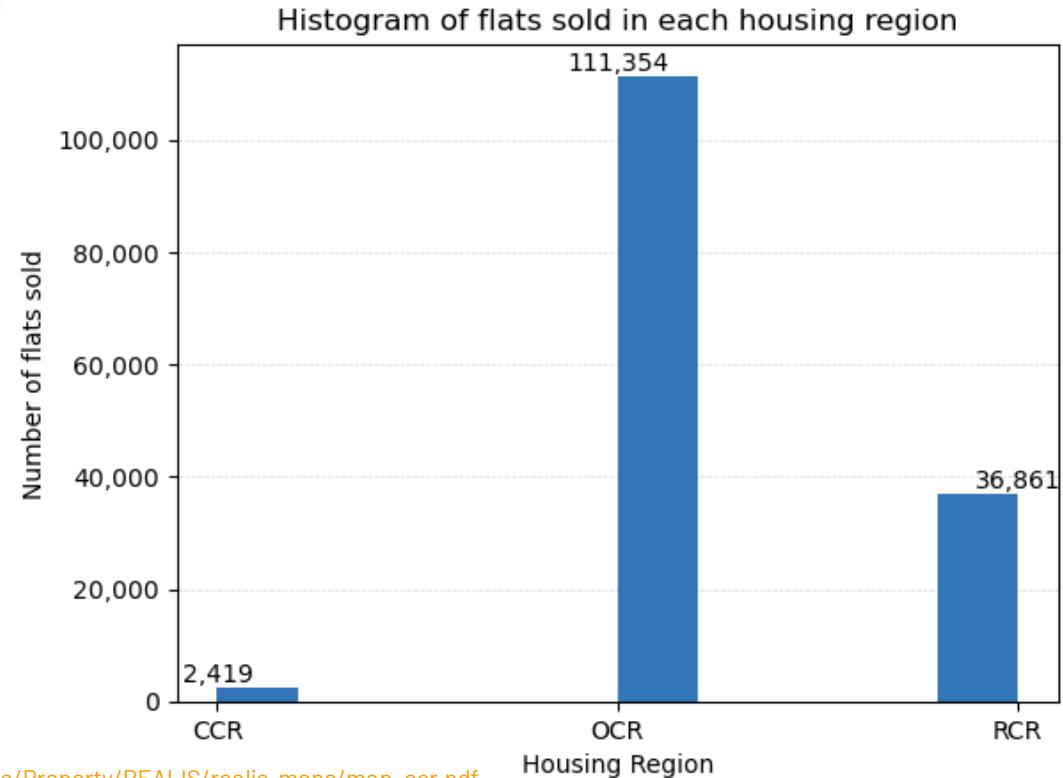
Legend\*:

CCR = Core Central Region

RCR = Rest of Central Region

OCR = Outside Central Region

**Most houses sold were in  
the OCR  
(Outside Central Region)**



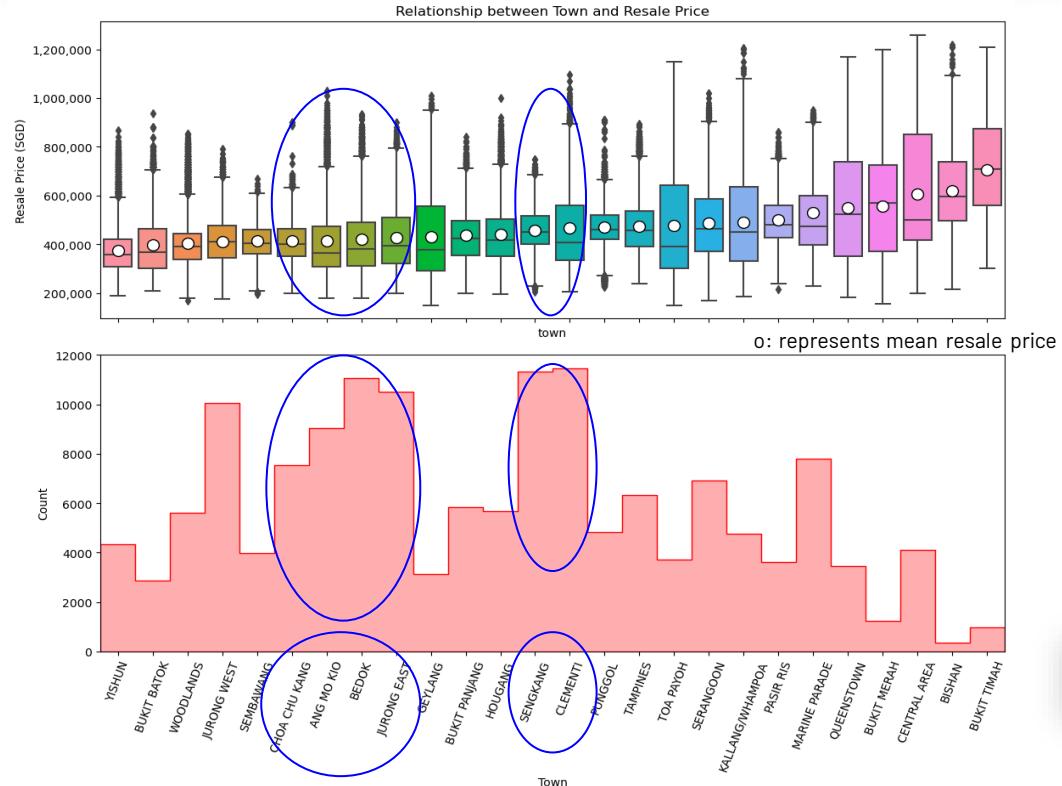
[https://www.ura.gov.sg/-/media/Corporate/Property/REALIS/realis-maps/map\\_ccr.pdf](https://www.ura.gov.sg/-/media/Corporate/Property/REALIS/realis-maps/map_ccr.pdf)

# Exploratory Data Analysis



Most houses sold from 2012 to 2021 are in the OCR region, not so much within central region

Towns with the more expensive resale prices had the lowest number of transactions



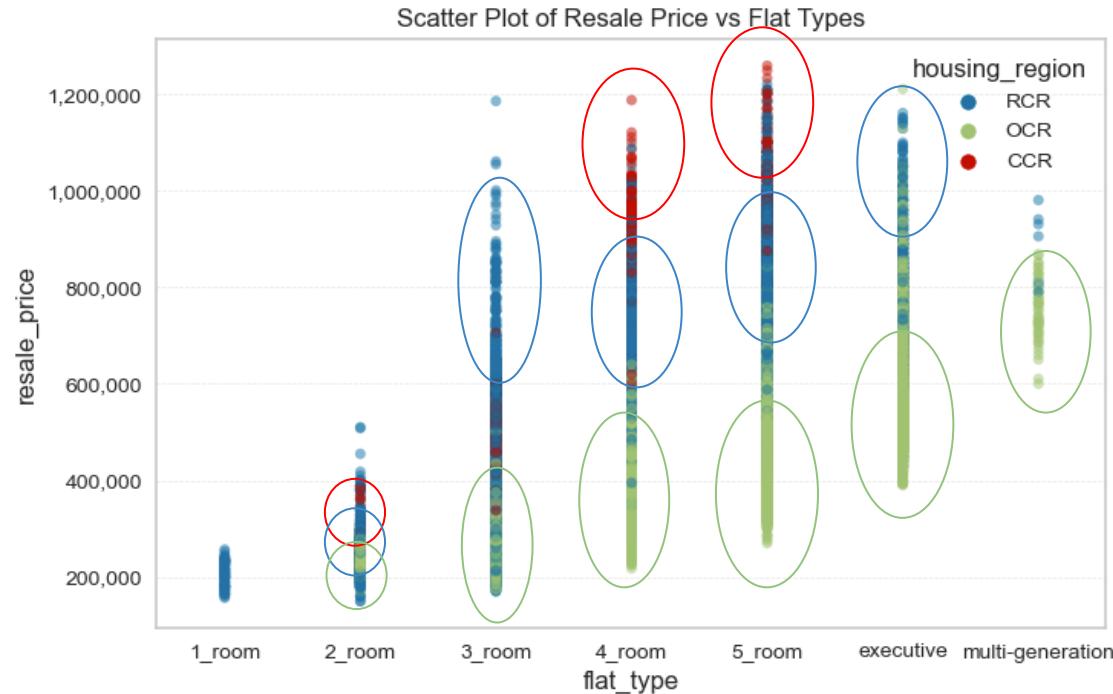
# Exploratory Data Analysis



## Legend:

- CCR = Core Central Region
- OCR = Outside Central Region
- RCR = Rest of Central Region

**CCR: Consistently more expensive**  
**RCR: Middleground**  
**OCR: Generally the least expensive**



# Exploratory Data Analysis



Low Floor:

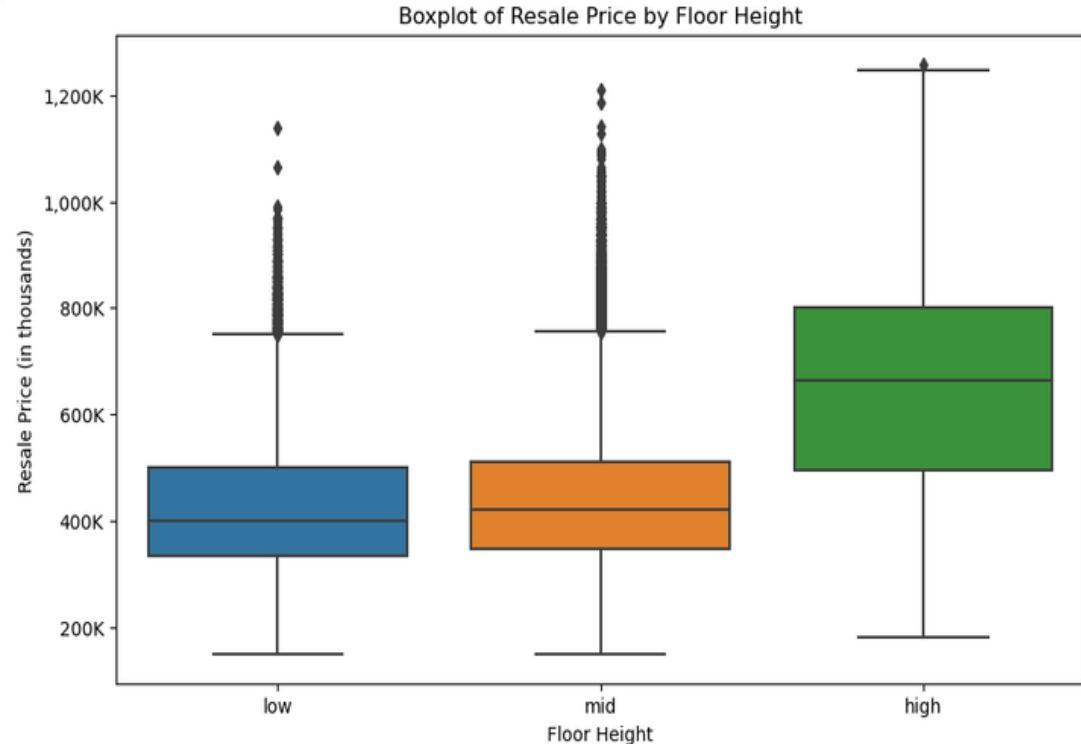
- Below 6 storey
- Less than  $\frac{2}{3}$  of low rise building
- Less than  $\frac{1}{3}$  of mid rise building

Mid Floor :

- Between storey 7 - 18
- Above  $\frac{1}{3}$  of mid rise building
- Above  $\frac{2}{3}$  of low rise building

High Floor:

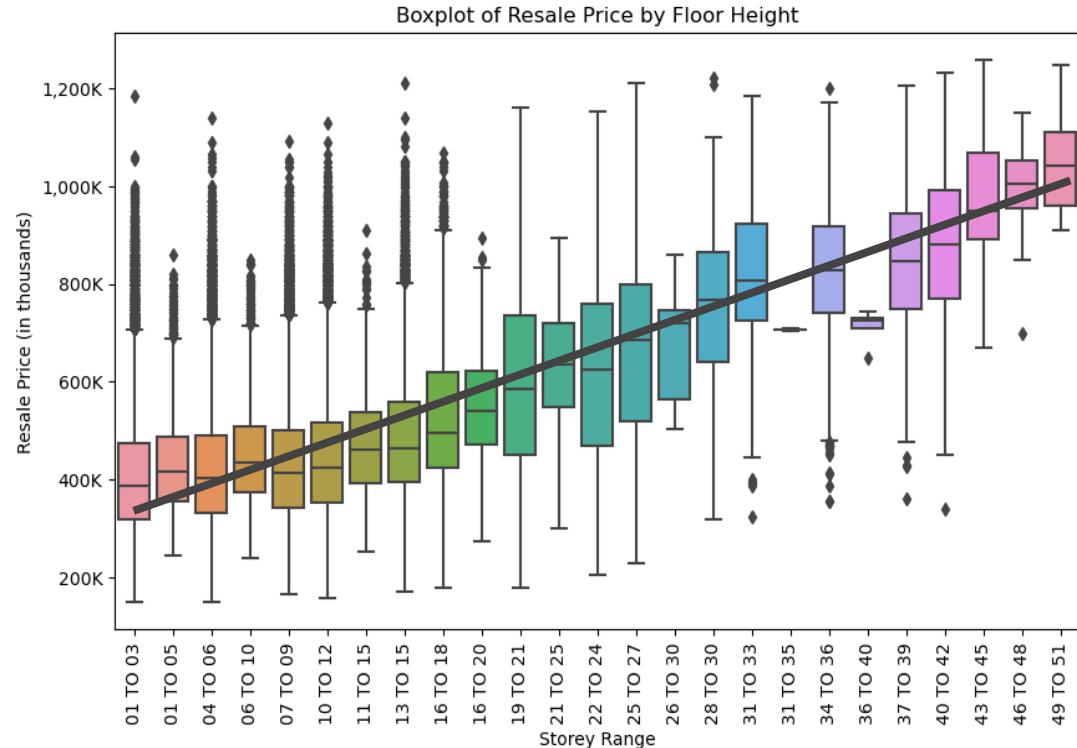
- Above storey 18



# Exploratory Data Analysis



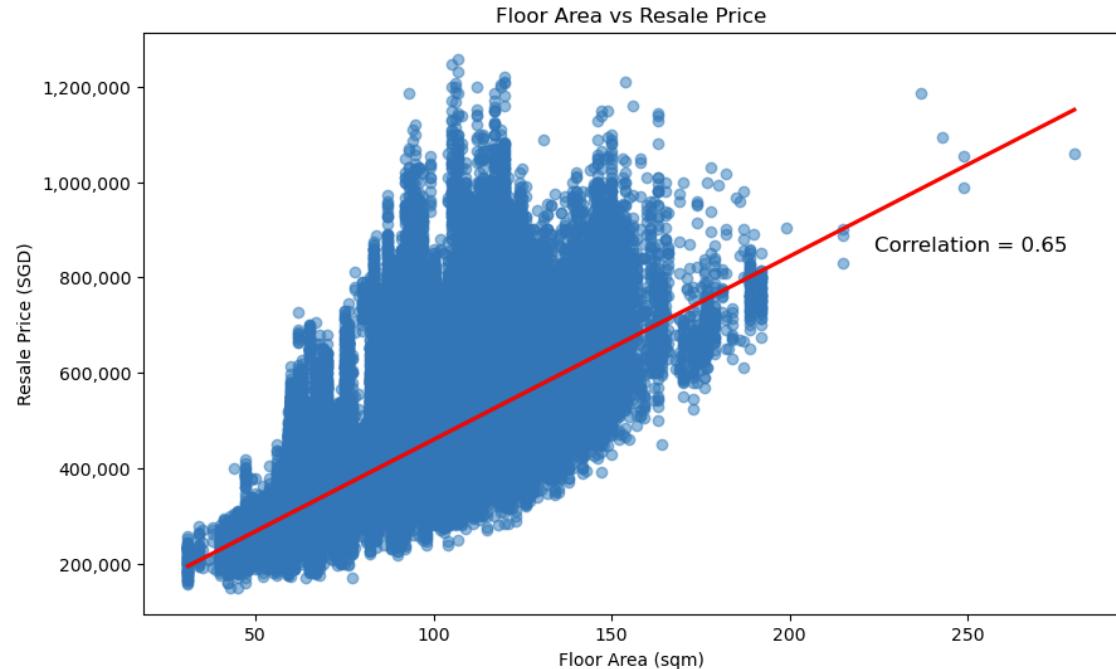
There's a general  
uptrend on the resale  
price as the storey range  
goes up



# Exploratory Data Analysis



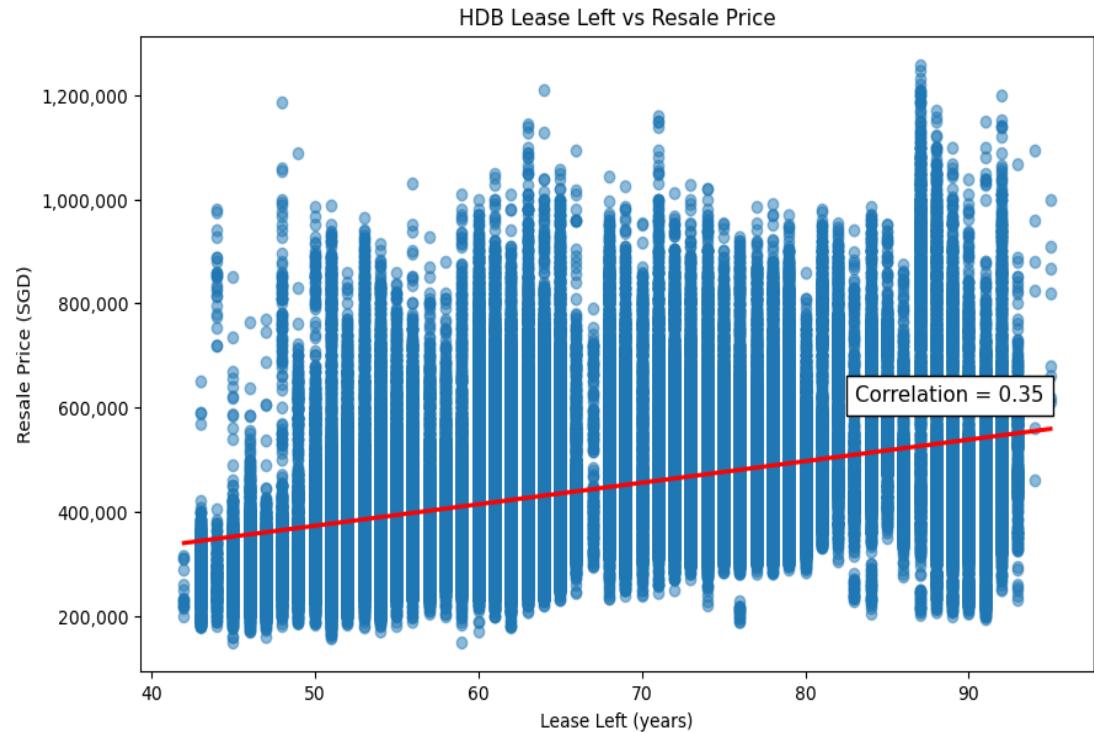
Floor Area of the HDB unit seems to have a strong positive correlation to the resale price of the HDB unit.



# Exploratory Data Analysis



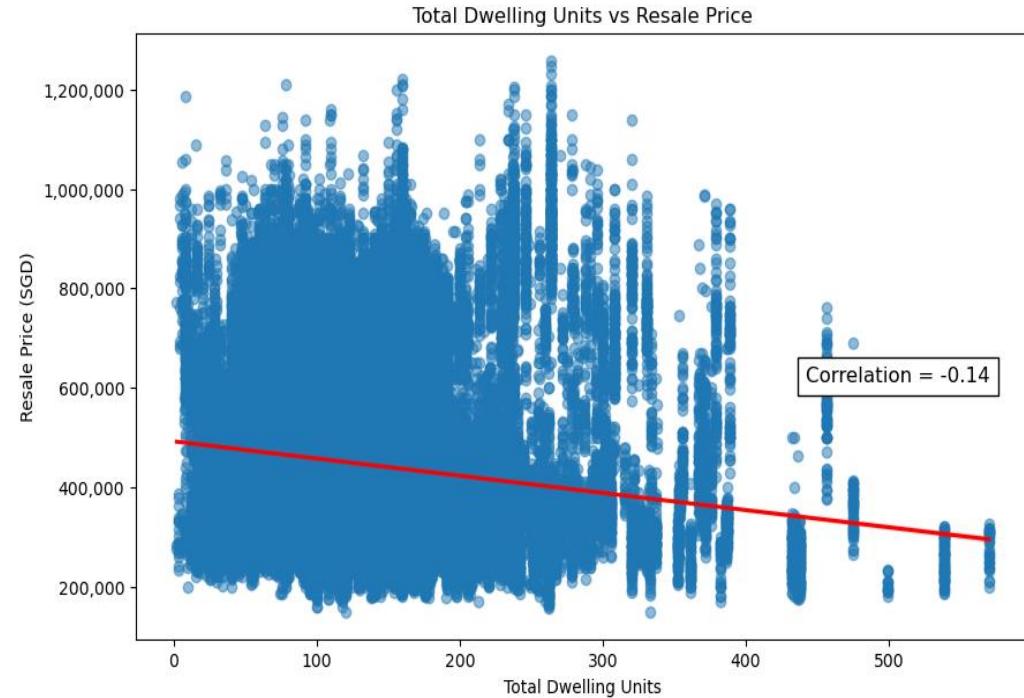
The more years left in the lease of the HDB unit is correlated to how high of a resale price the HDB unit is able to fetch.



# Exploratory Data Analysis



Total dwelling units in a HDB block is not observed to have a strong correlation with the resale price.

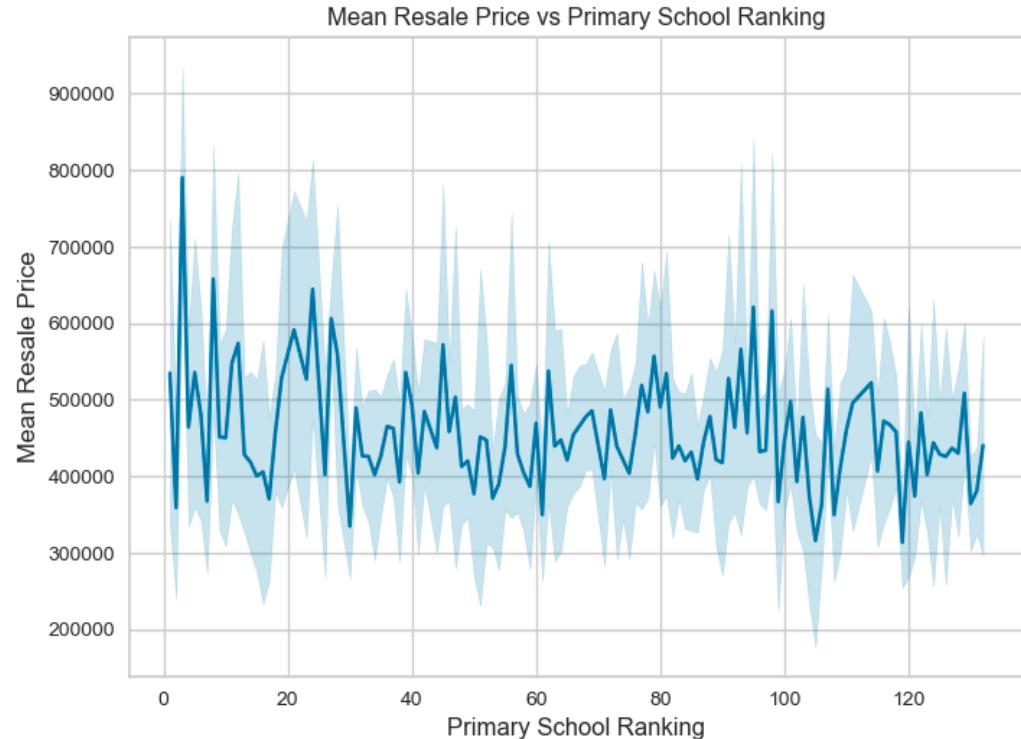


# Exploratory Data Analysis



Nearby primary school ranking does not seem to have an obvious influence over the resale price of a HDB unit.

This observation is consistent with the findings of other sources.

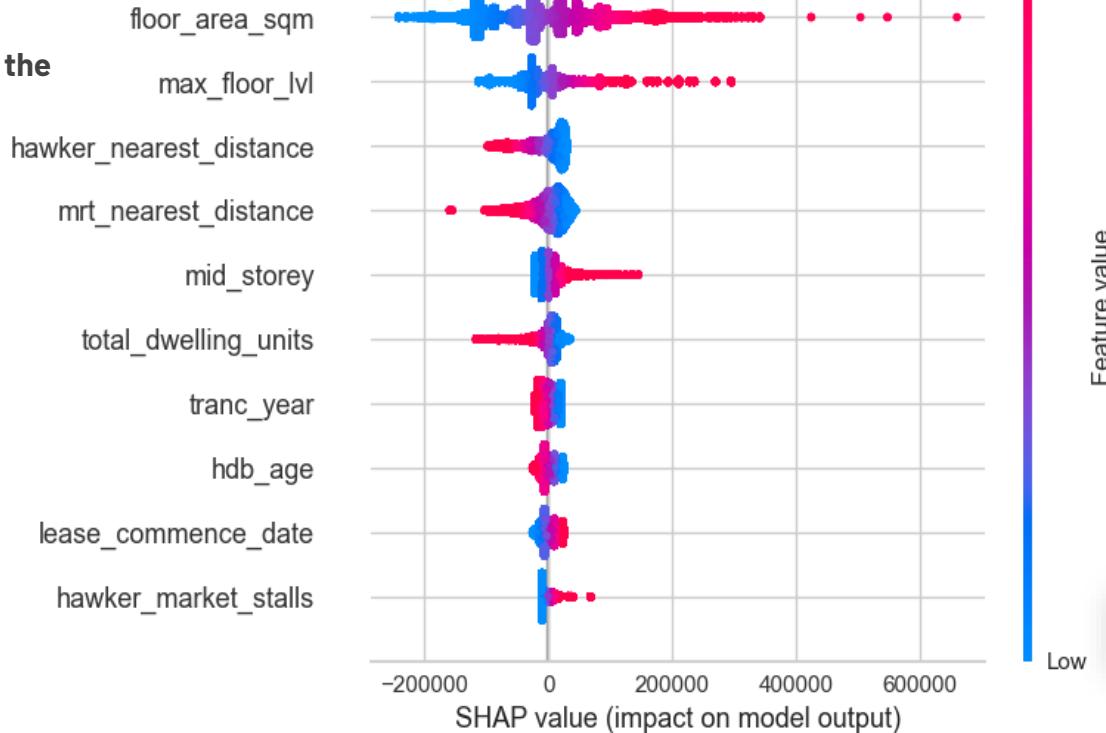


Source: <https://dollarsandsense.sg/hdb-property-prices-near-popular-primary-schools-really-cost/>

# Exploratory Data Analysis



**Features that were found to have the greatest impact on the model output .**





# LINEAR REGRESSION MODEL

Building a “linear guideline” for making predictions of future prices of HDB units.

# Regression Model Results



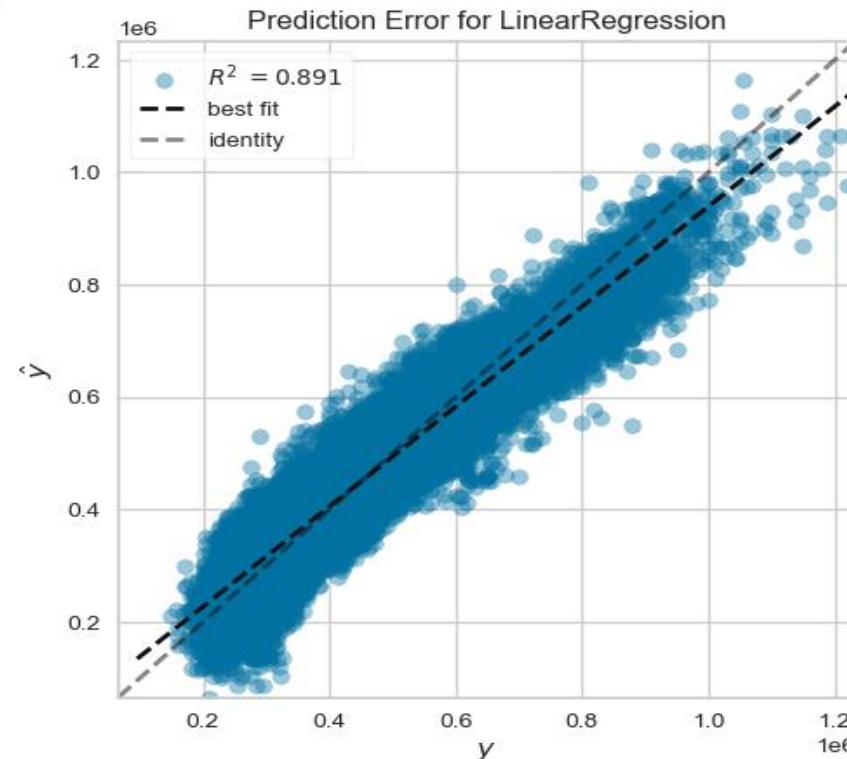
Model 1 performed better on unseen data set

Model	R <sup>2</sup> (train)	R <sup>2</sup> (test)	RMSE	RMSE (Unseen data, kaggle)
<b>Baseline</b>	0.8638	0.8610	52,965.62	-
<b>1</b>	0.8911	0.8904	47,162.25	47,149.78
<b>2</b>	0.9147	0.9137	41,837.78	54,155.30
<b>3</b>	0.8927	0.8920	46,806.20	55,781.83

# Linear Regression Model



**Prediction error plot shows the actual targets from the dataset against the predicted values generated by our model**



# Linear Regression: Feature Selection



Model	Feature Selection Description
<b>Baseline</b>	<ul style="list-style-type: none"><li>• Baseline model runs with all numeric features</li><li>• Used as a baseline to evaluate model performance</li></ul>
<b>1</b>	<ul style="list-style-type: none"><li>• Feature selection based on domain knowledge</li><li>• Elements that are known to affect housing prices</li></ul>
<b>2</b>	<ul style="list-style-type: none"><li>• The features selection are based on features correlation</li><li>• Feature engineering of region against flat types</li><li>• Popularity ranking of primary schools</li><li>• Availability of amenities</li></ul>
<b>3</b>	<ul style="list-style-type: none"><li>• Feature selection based on model 1 features and</li><li>• Feature importance from previous models.</li></ul>



# PRODUCT DEMONSTRATION

Live demonstration of your property valuation



# **RECOMMENDATION & KEY INSIGHTS**

So what's the gist of it?

# Recommendations



## Buyer

- Know your available options given your budget
  - Buyers should be able to make an informed decision that fits their budget.
- Prioritize and personalize your wants
  - Buyers are recommended to straighten out their priority and decide on the factors that fits their needs most.
- Quality home with comfortable price

## Seller

- Appraise your property value based on market valuation
  - Sellers are recommended to at least have some understanding of the market resale price of their respective units.
- Pivot your selling strategies
  - Sellers are recommended to be ready to switch up with their selling strategies at any given time due to market volatility.
- Match your property's unique selling points to the right buyers

# 1 Minute

To have an estimated price for your dream house

# 89.1%

Prediction accuracy

# 11 X-Factors

Focus on the factors that matters



# KEY LIMITATIONS

If only we had more time and resource.

# Key Limitations



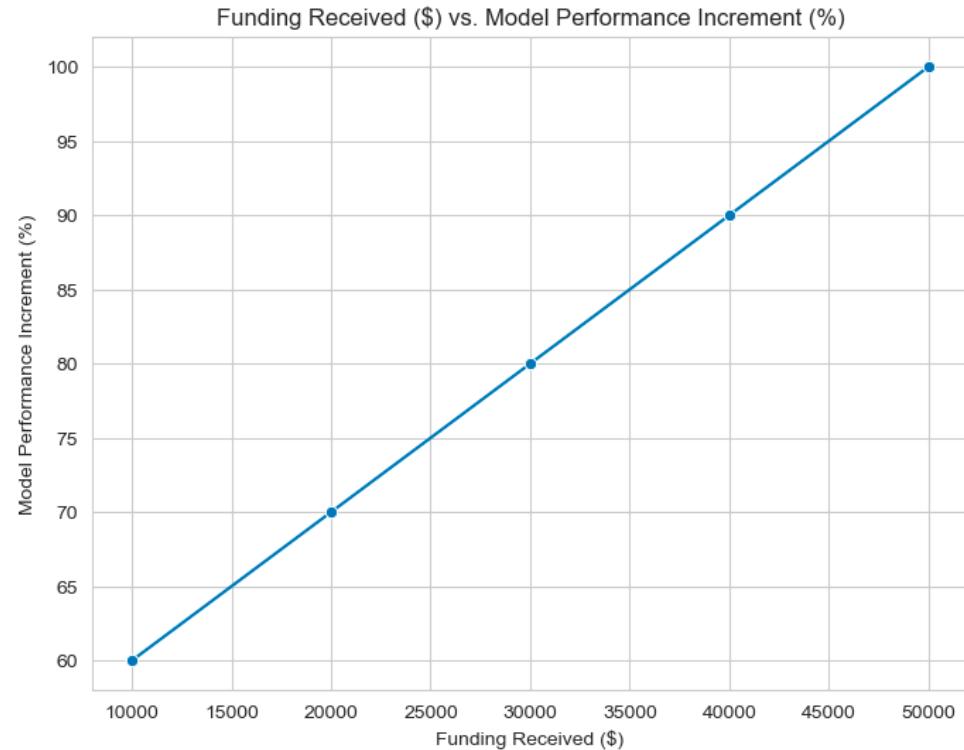
- Limitation of data
  - Data collected is only up to 2021
- Modelling is limited to only Linear Regression
  - Opportunity to utilize more sophisticated model in the future for better prediction.
- Lack of info on the existing condition of the sold units
- Collinearity does not imply causation

# Key Limitations (Joke)



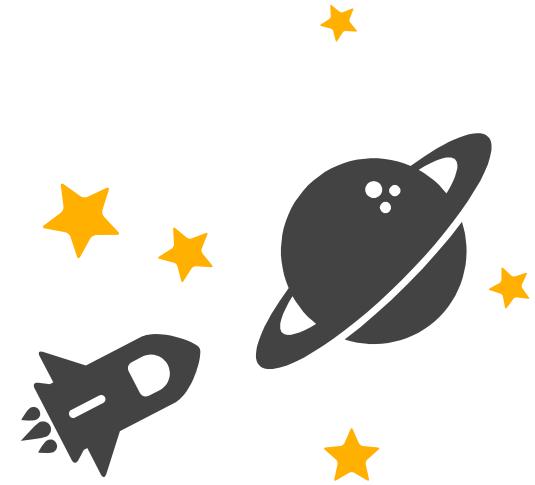
## Funding

Severe lack of funding is making it challenging to carry out further refinement works to increase the model performance.



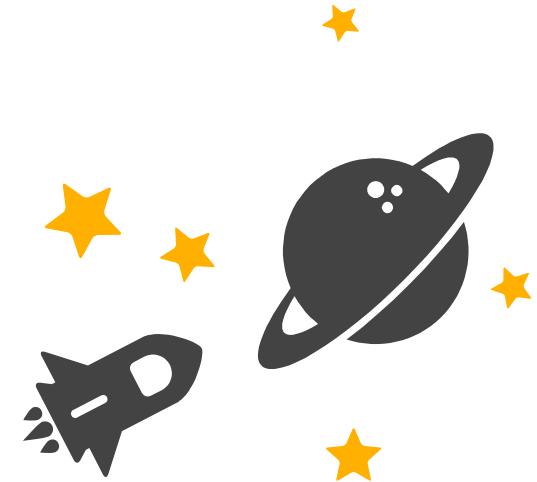
# Your House, Your Future

Make your real estate plans with  
technology of future





# THANK YOU!





# **BACK-UP SLIDES / ANNEX**

Extra cheese

# Linear Regression Model



Model 1 Features:

Target (y-axis): resale\_price

X(axis):

1. `town` (cat)
2. `storey\_range` (cat)
3. `full\_flat\_type` (cat)
4. `pri\_sch\_name` (cat)
5. floor\_area\_sqm
6. lease\_commence\_date
7. mrt\_nearest\_distance
8. hawker\_nearest\_distance
9. mall\_nearest\_distance
10. pri\_sch\_nearest\_distance
11. sec\_sch\_nearest\_dist

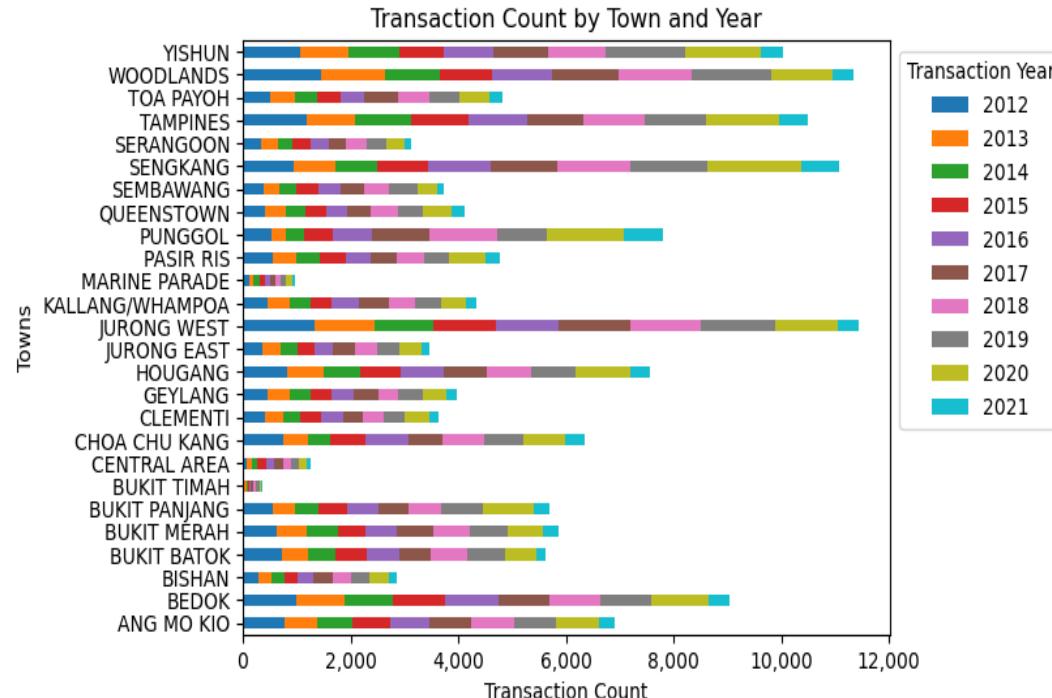
## Features

Train.csv features: 78

Dropped: id, price\_per\_sqft, floor\_area\_sqft, resale\_price (which is target)

Added: pop\_ranking, pop\_ranking\_2cat, housing\_region

# Exploratory Data Analysis

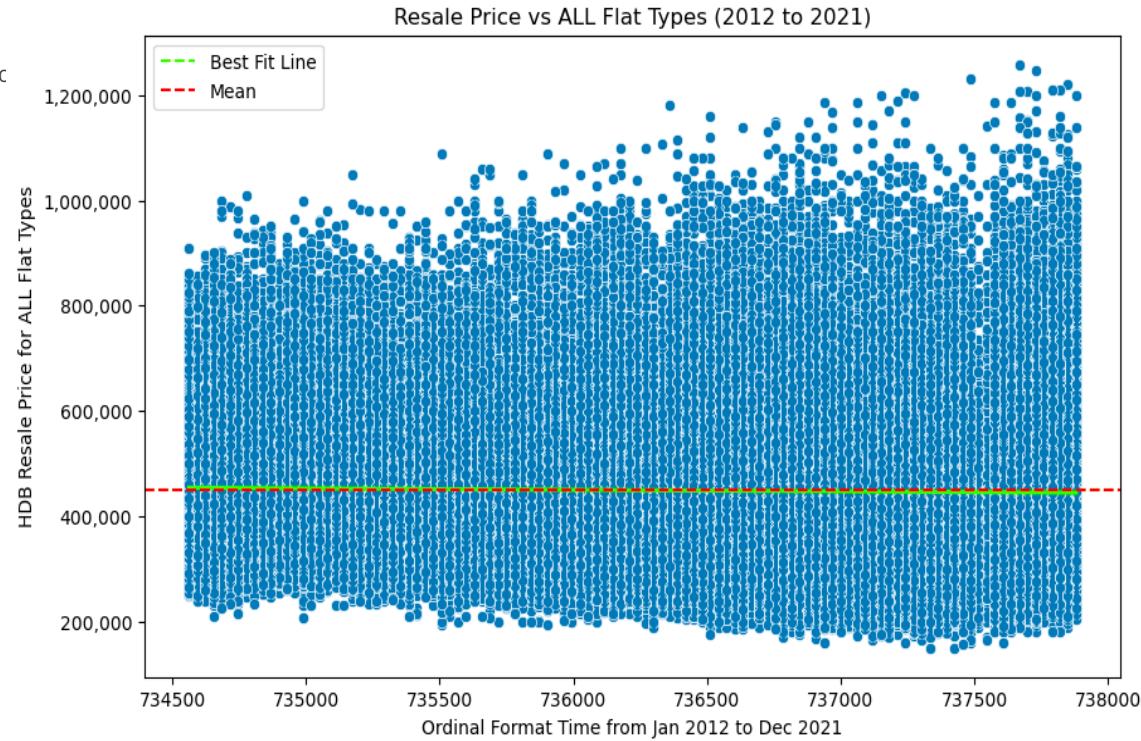


# Exploratory Data Analysis



Note:

Progression of years is found to not have much effect on the resale prices of HDB units

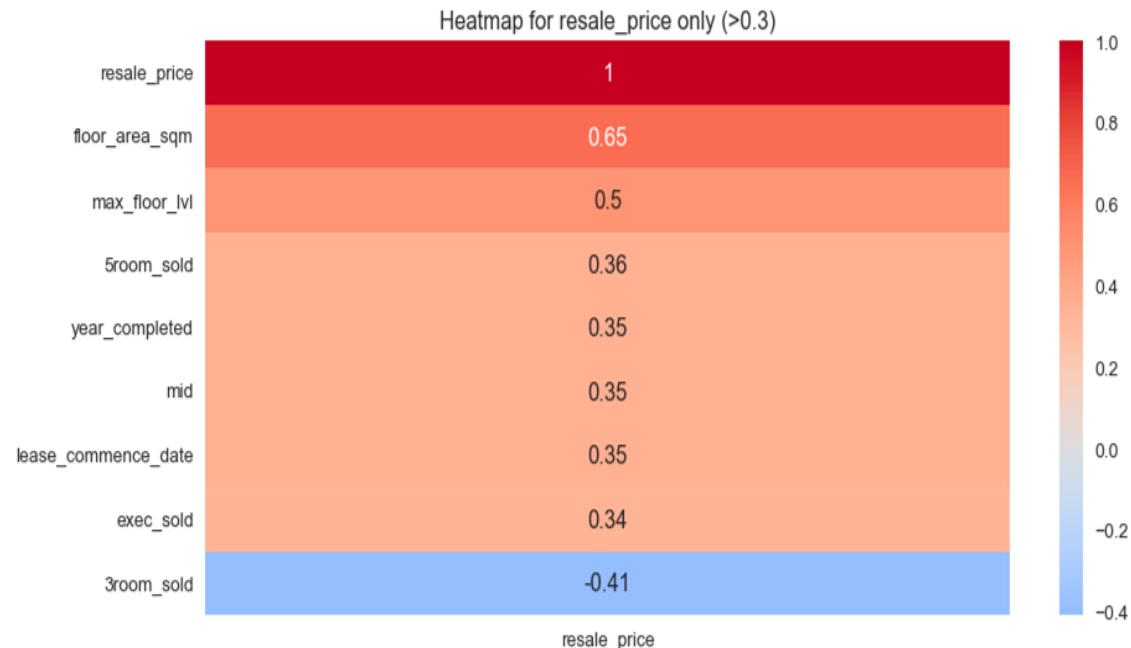


# Exploratory Data Analysis

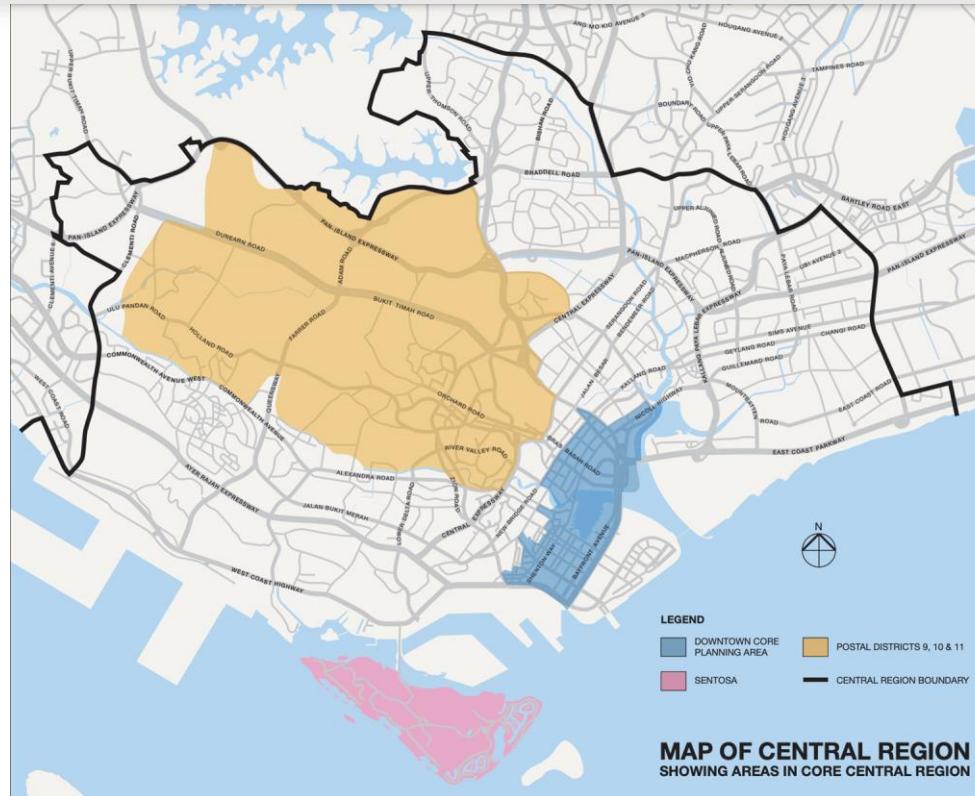


Note:

Heatmap of correlation coefficient of features to resale\_price.



# Housing Regions: CCR, RCR, OCR



URA Map from:

[https://www.ura.gov.sg/-/media/Corporate/Property/REALIS/realis-maps/map\\_ccr.pdf](https://www.ura.gov.sg/-/media/Corporate/Property/REALIS/realis-maps/map_ccr.pdf)

## Postal Sector



- <https://www.mingproperty.sg/singapore-district-code/>

# Regression Model Results



Linear Regression Model	R <sup>2</sup> (train)	R <sup>2</sup> (test)	RMSE	RMSE (Unseen data, kaggle)
Baseline	<b>0.8638</b>	<b>0.8617</b>	<b>52,965.62</b>	
1	<b>0.8911</b>	<b>0.8904</b>	<b>47,162.25</b>	<b>47,149.78</b>
2	<b>0.9147</b>	<b>0.9137</b>	<b>41,837.78</b>	<b>47,509.99</b>
3	<b>0.8917</b>	<b>0.8920</b>	<b>46,806.20</b>	<b>55,781.83</b>

# Regression Model Results



Ridge Regression Model	R <sup>2</sup> (train)	R <sup>2</sup> (test)	RMSE	RMSE (Unseen data, kaggle)
Baseline	<b>0.8638</b>	<b>0.8617</b>	<b>52,965.61</b>	
1	<b>0.8911</b>	<b>0.8903</b>	<b>47,163.40</b>	<b>47,149.78</b>
2	<b>0.9146</b>	<b>0.9136</b>	<b>41,854.66</b>	<b>47,509.99</b>
3	<b>0.8927</b>	<b>0.8920</b>	<b>46,805.96</b>	<b>55,781.83</b>

# Regression Model Results



Lasso Regression Model	R <sup>2</sup> (train)	R <sup>2</sup> (test)	RMSE	RMSE (Unseen data, kaggle)
Baseline	<b>0.8631</b>	<b>0.8610</b>	<b>53,105.47</b>	
1	<b>0.5761</b>	<b>0.5751</b>	<b>92,861.14</b>	<b>47,149.78</b>
2	<b>0.7071</b>	<b>0.7046</b>	<b>77,422.44</b>	<b>47,509.99</b>
3	<b>0.8920</b>	<b>0.8914</b>	<b>46,927.46</b>	<b>55,781.83</b>

# Linear Regression Model



Note:  
Model 2

