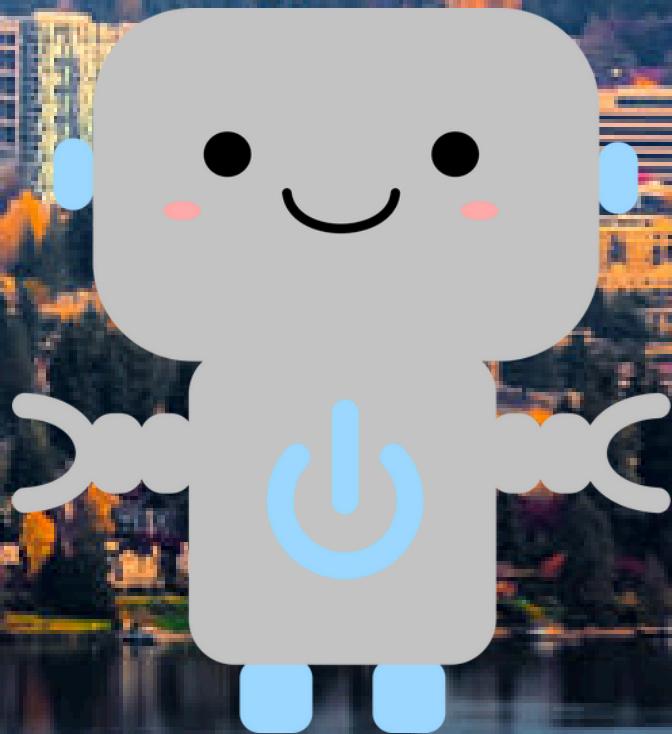


# SC1015: Mini Project

Team 5:

Jumana, Aung, Yichi



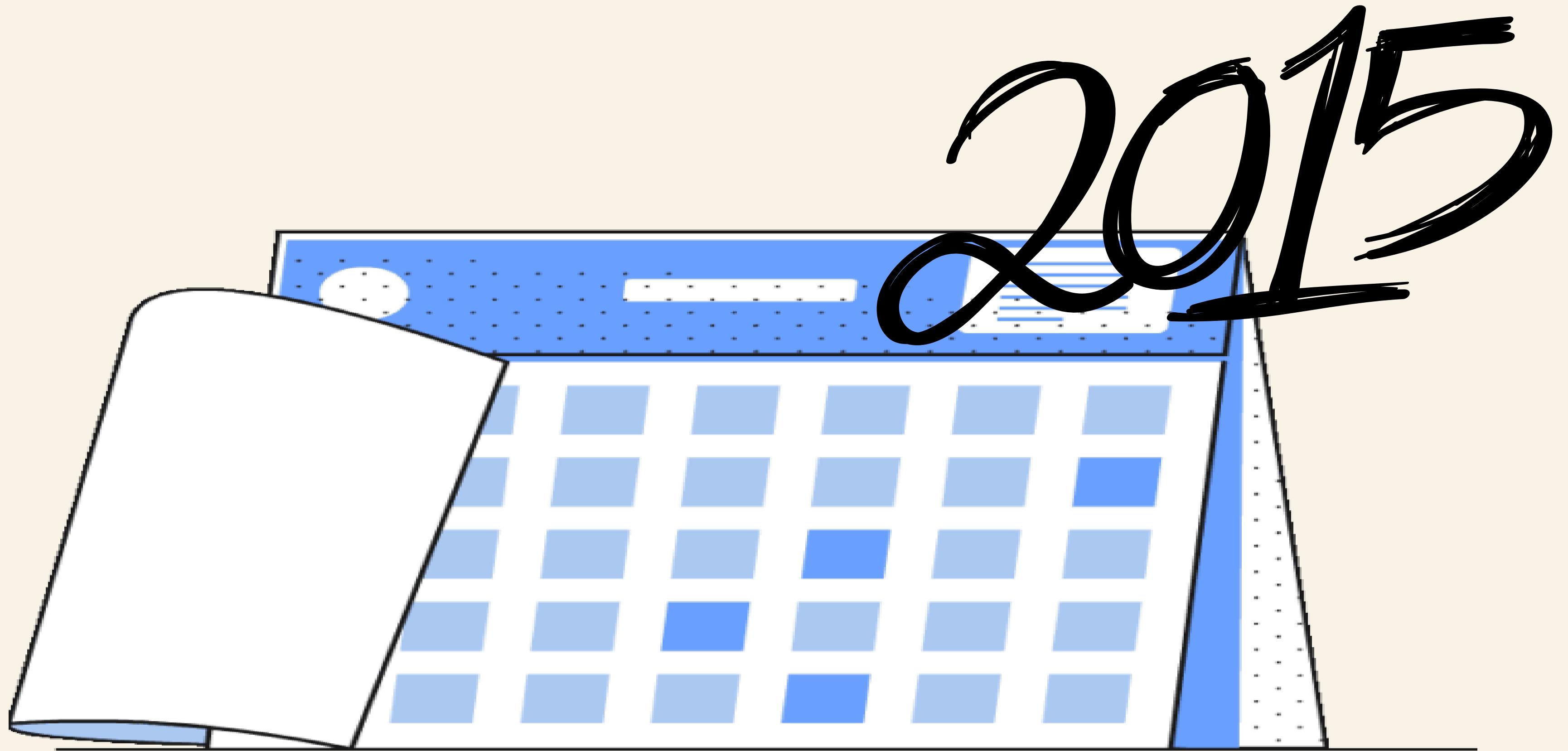
01 - Problem formulation

02 - Exploratory Data Analysis

03 - Machine Learning

04 - Insights & Recommendations



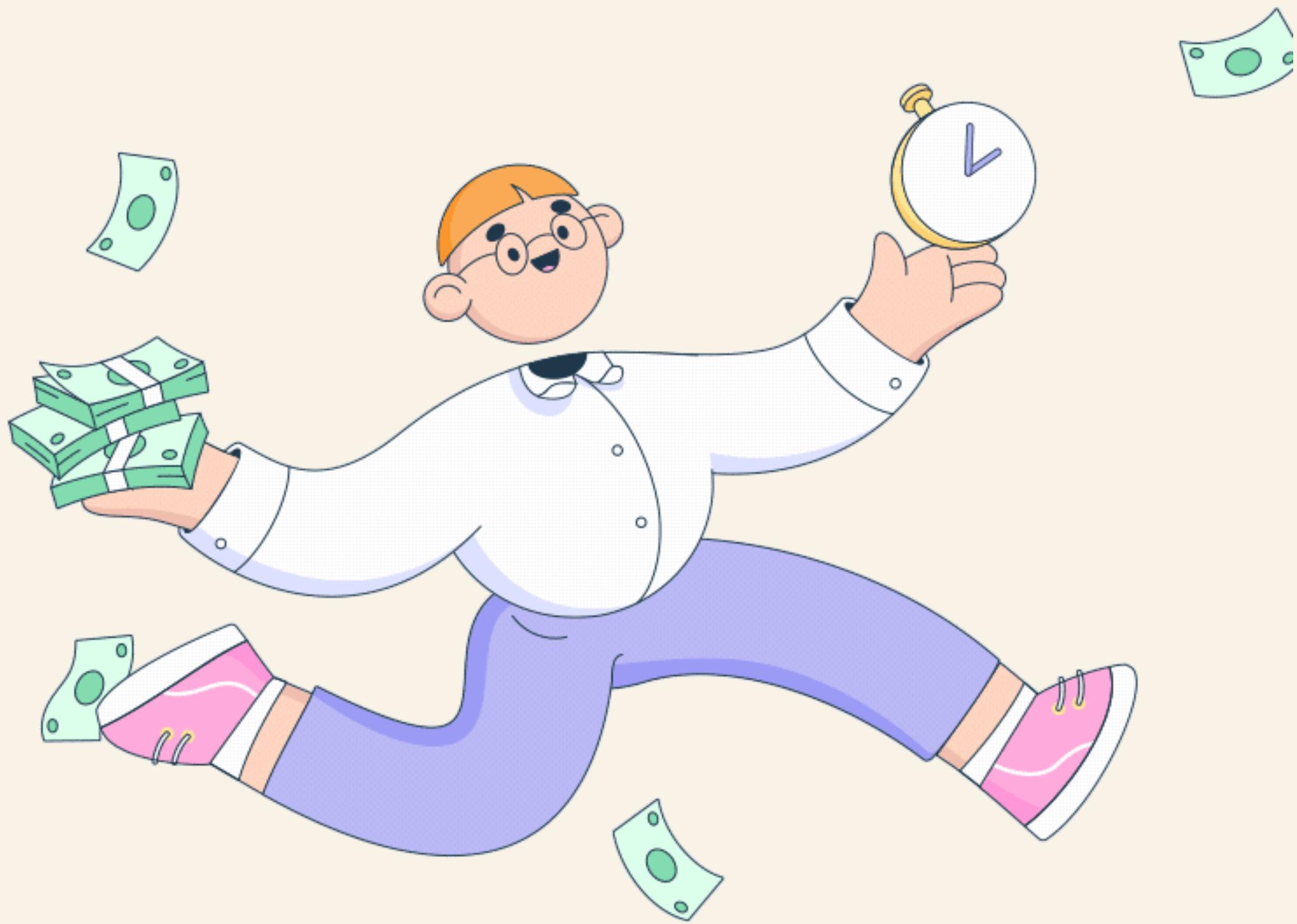


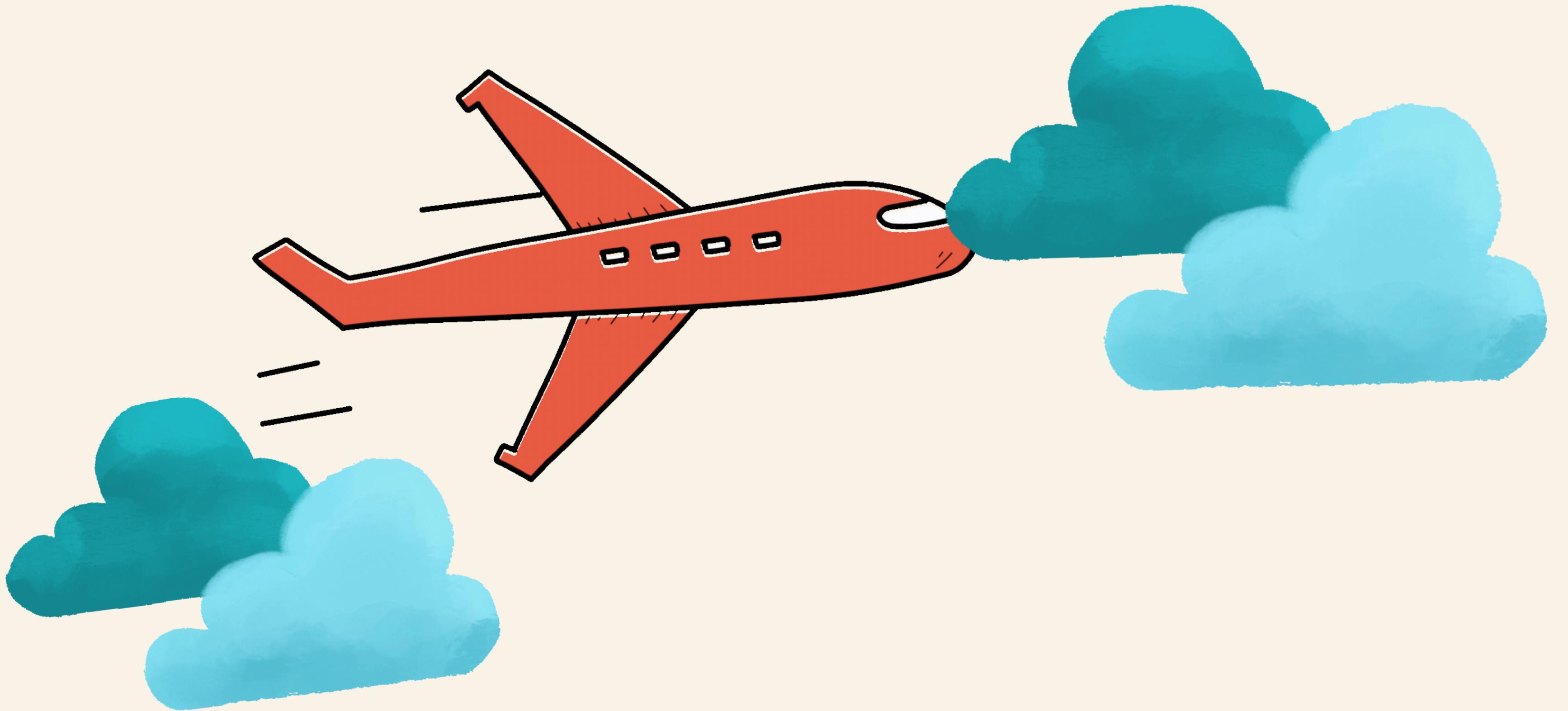


Lucas

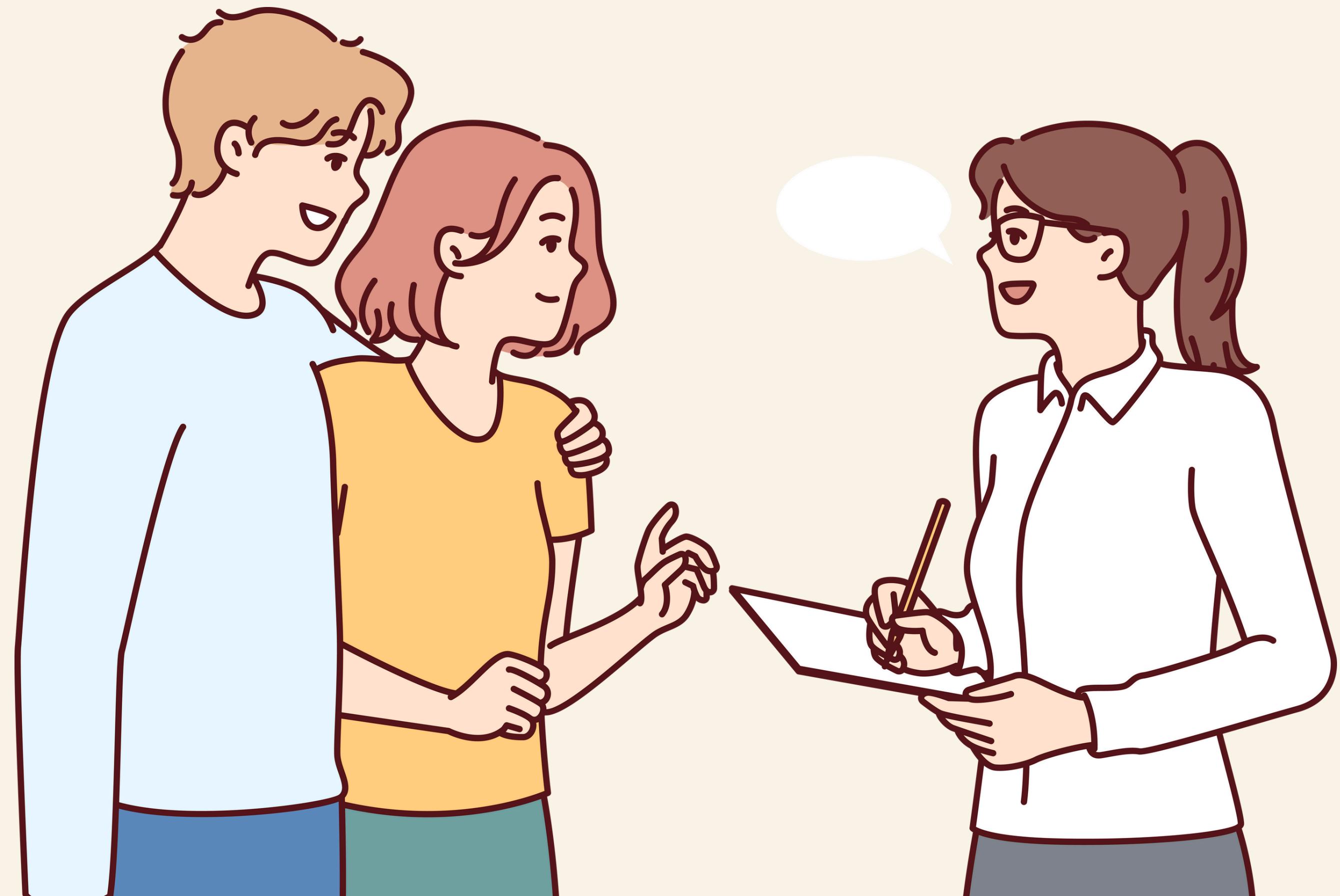


Lucas









# Problem formulation

# Problem formulation



# Problem formulation



# kaggle

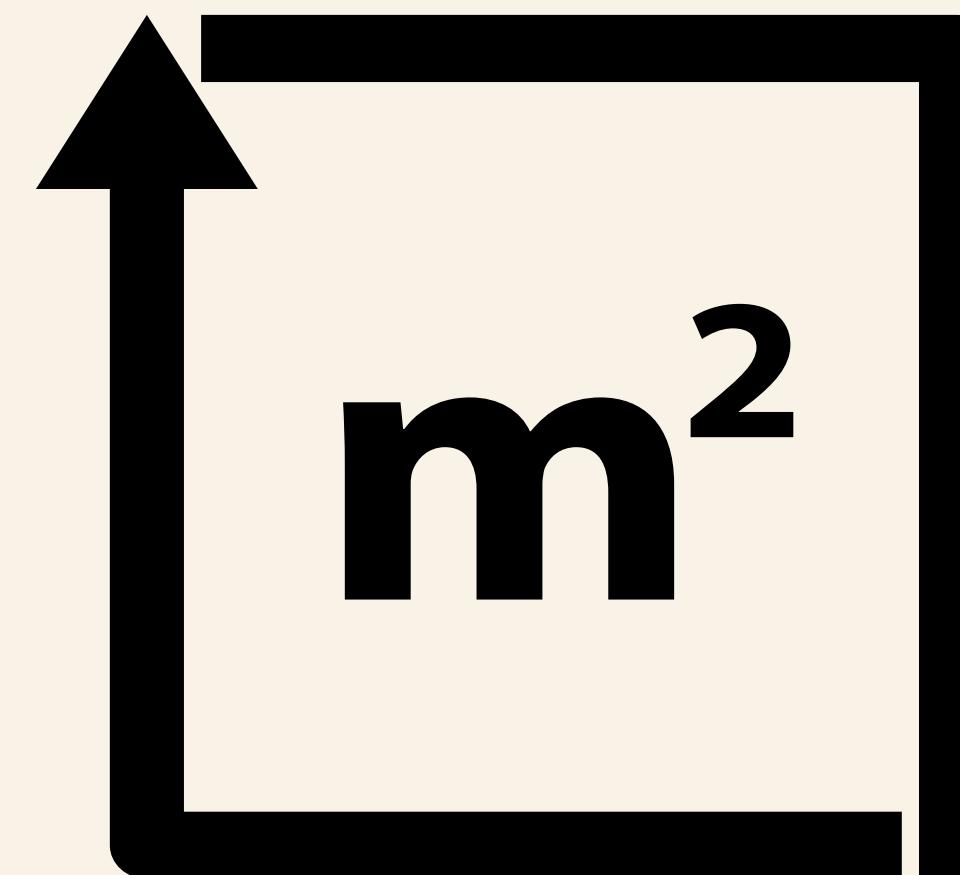
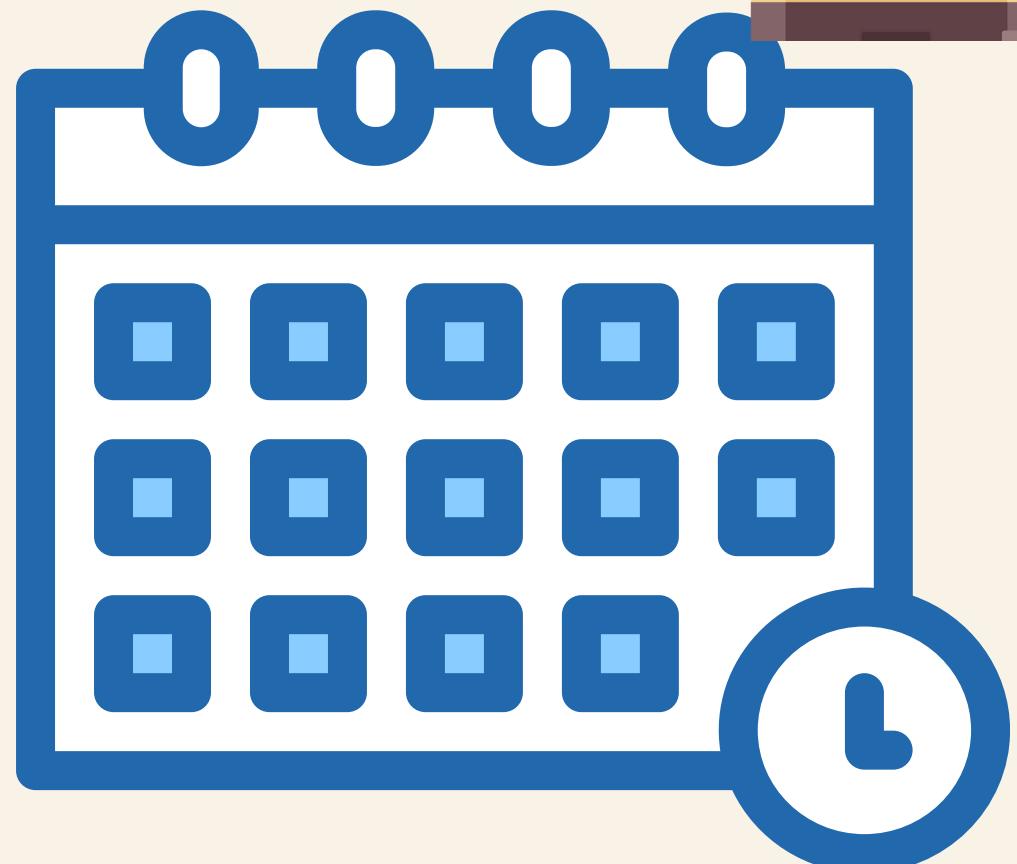
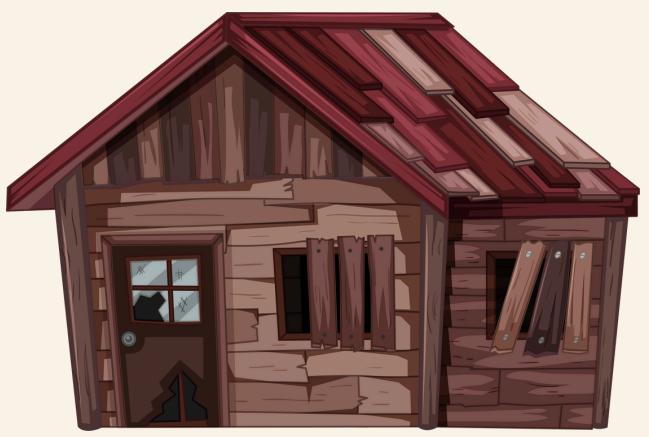


Variable	Description
id	Unique ID per house sale
date	Date of the house sale
price	Price of house sale in currency of USD
bedrooms	Number of bedrooms
bathrooms	Number of Bathrooms, where 0.5 represents a bathroom with a toilet but with no shower
sqft_living	Square footage of the apartments interior living space
sqft_lot	Square footage of the land space
floors	Number of floors
waterfront	An index to indicate if the house was overlooking the waterfront or not. 0 represents no waterfront, 1 represents with waterfront.
view	An index from 0 to 4 of how good the view of the property was. 0 represents no good view, 4 represents very good view.
condition	An index from 1 to 5 on the condition of the house. 1 represents poorer condition, and 5 represents superb condition.
grade	An index from 1 to 13. 1 to 3 falls short of building construction and design, 7 has an average level of construction and design, and 11 to 13 have higher quality level of construction and design.
sqft_above	The square footage of the interior housing space that is above the ground level
sqft_basement	The square footage of the interior housing space that is below the ground level
yr_built	The year of house built
yr_renovated	The year of the house's last renovation
zipcode	The zipcode is the postal code to indicate the area the house is in
lat	Latitude
long	Longitude
sqft_living15	The average square footage of interior housing living space for the nearest 15 neighboring houses
sqft_lot 15	The average square footage of land space for the nearest 15 neighboring houses

# Investopedia



# Investopedia



# Questions to ask

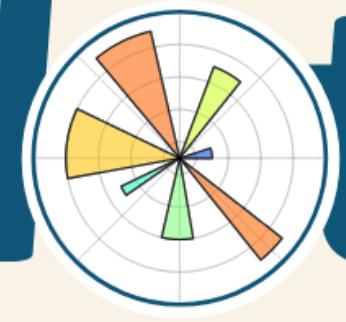
- What factors affect the changes in house sale price?
- How do house sale prices vary in King County?
- What might be the pattern of houses sold in different parts of King County?
- Is there any specific time period where houses are more expensive/cheaper?

# Exploratory Data Analysis



seaborn

matplotlib



plotly

## Trend of Average House Sale Price in King County Based on Year It Was Built (1900-2015)

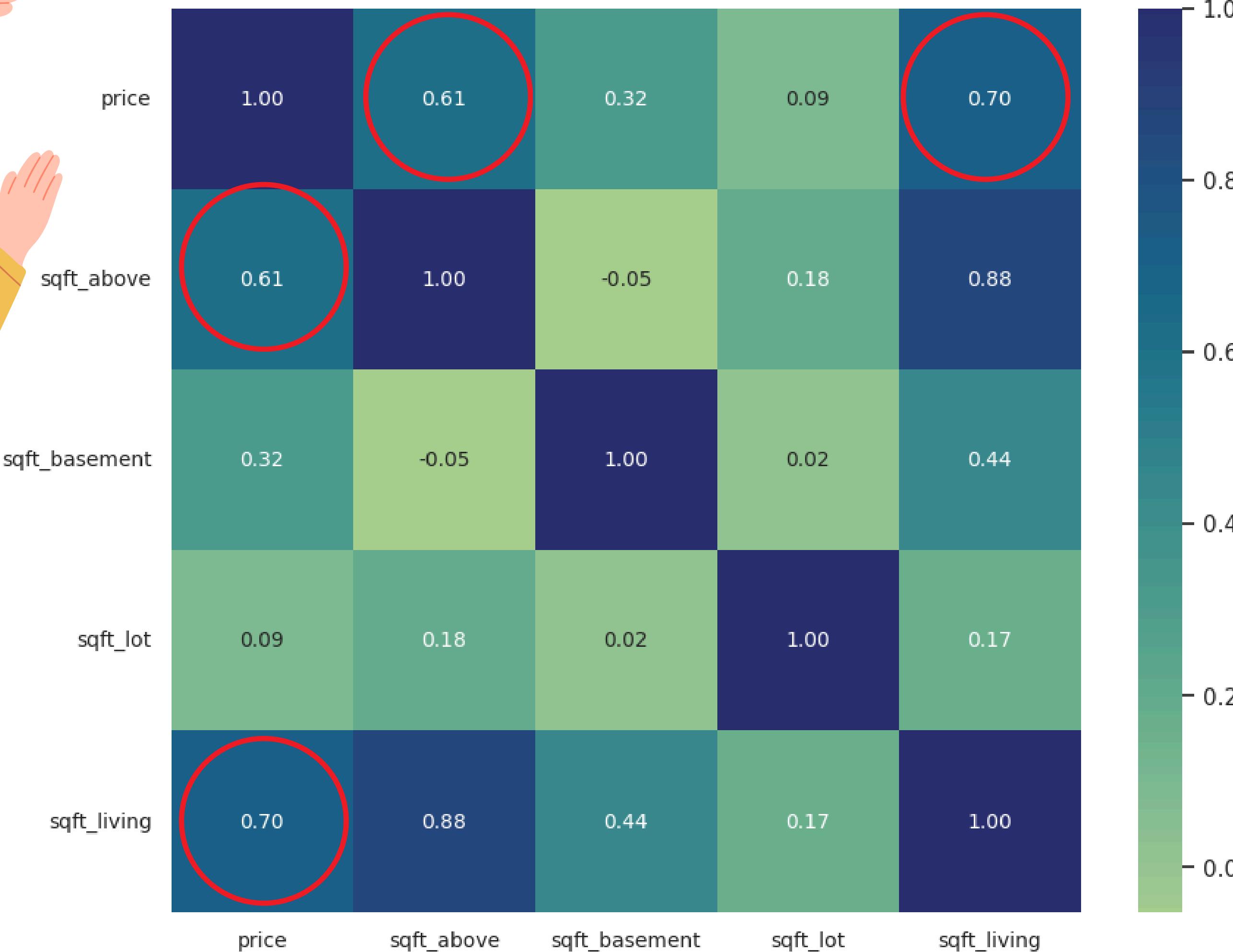


## Trend of Average House Sale Price in King County Based on Month It Was Sold

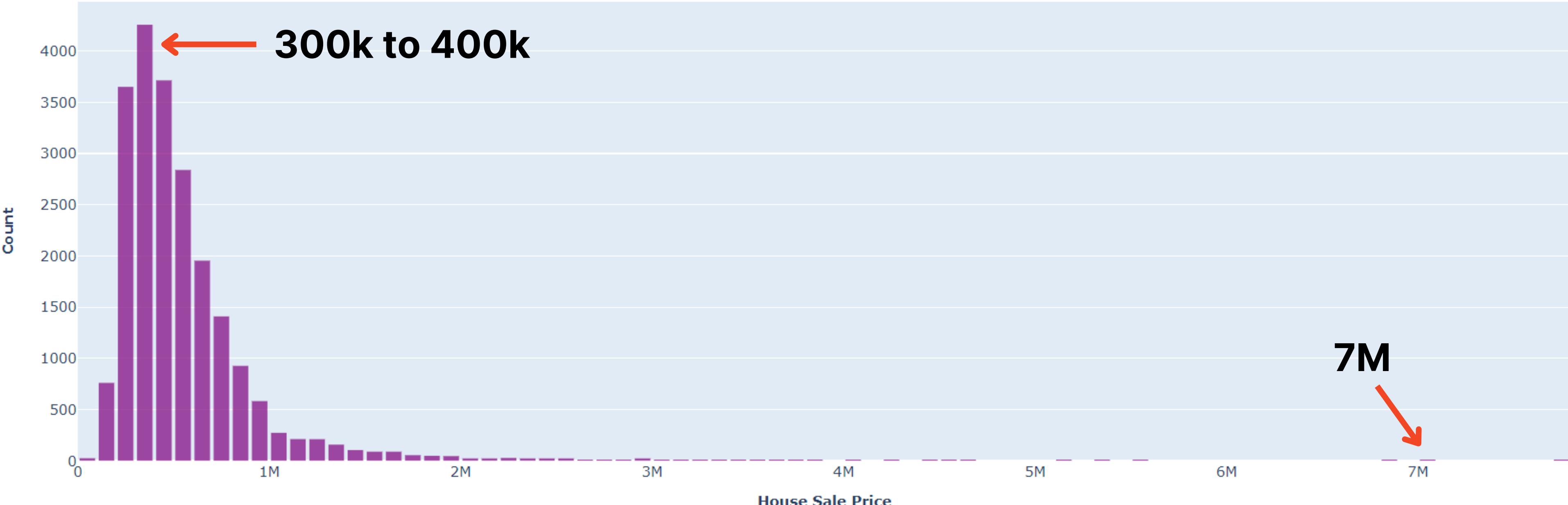




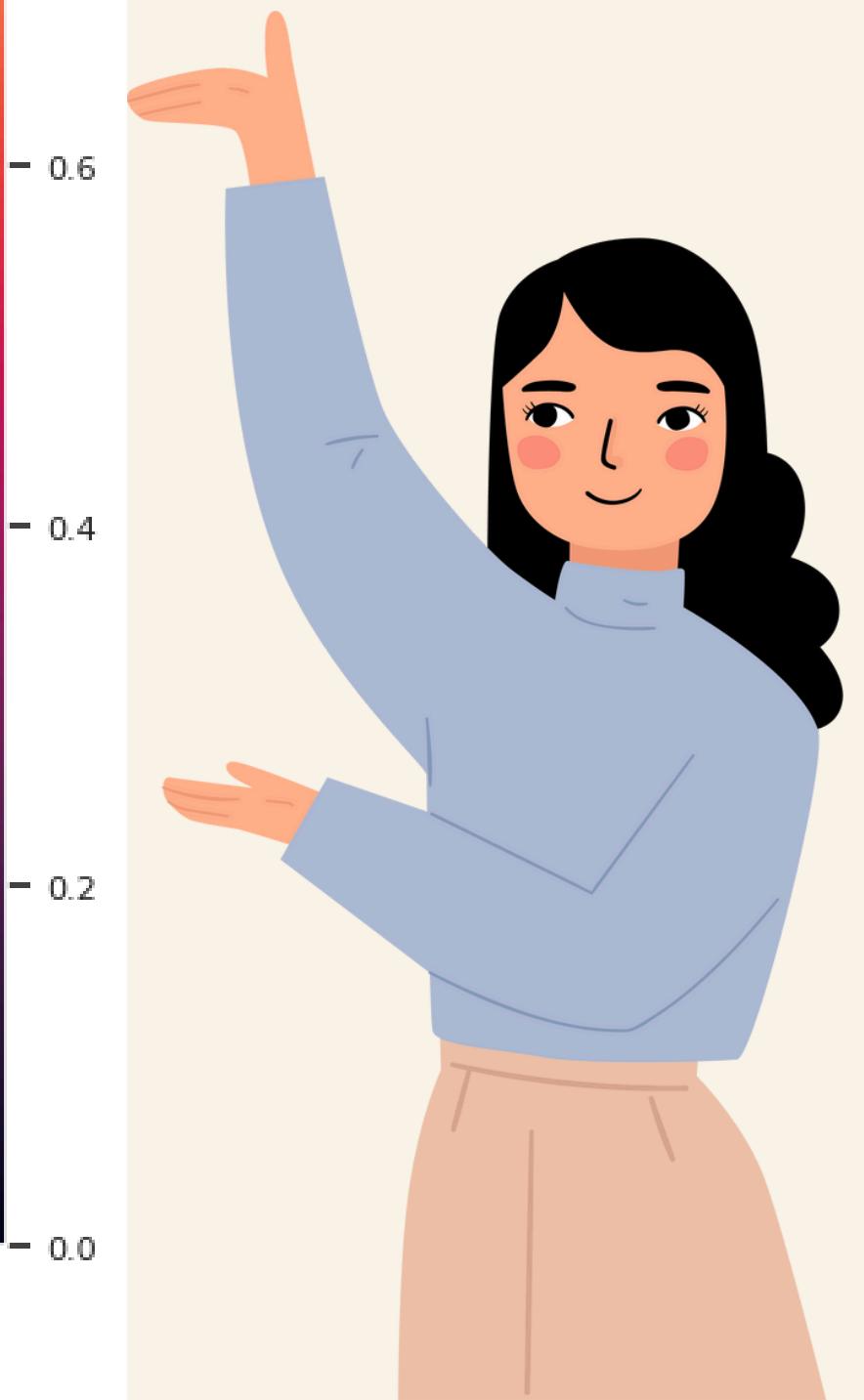
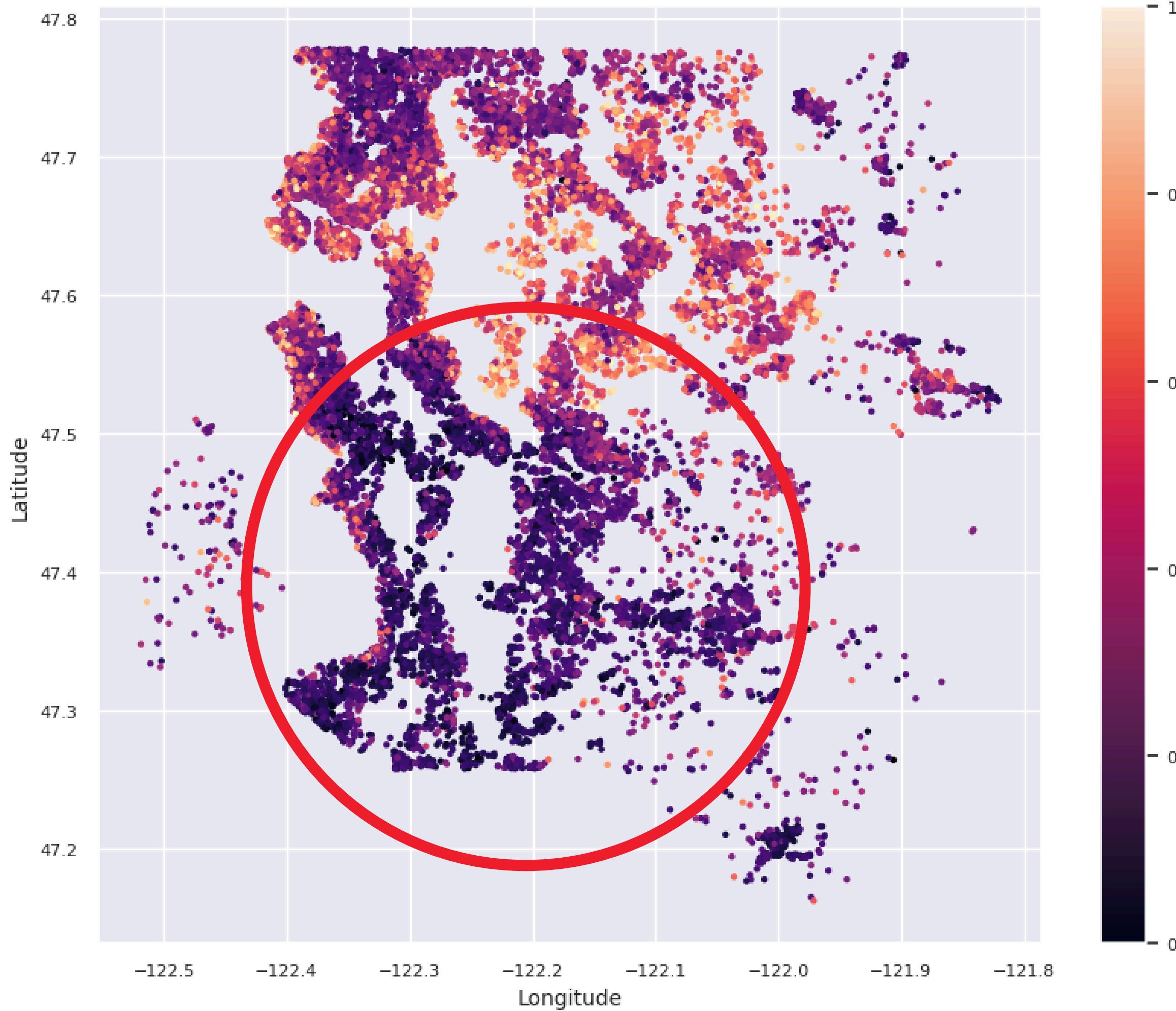
### Correlation of Sale Price And Different Measures of Houses in King County



## Distribution of House Sale Price in King County



# General Changes in House Sale Price Across Different Parts of King County



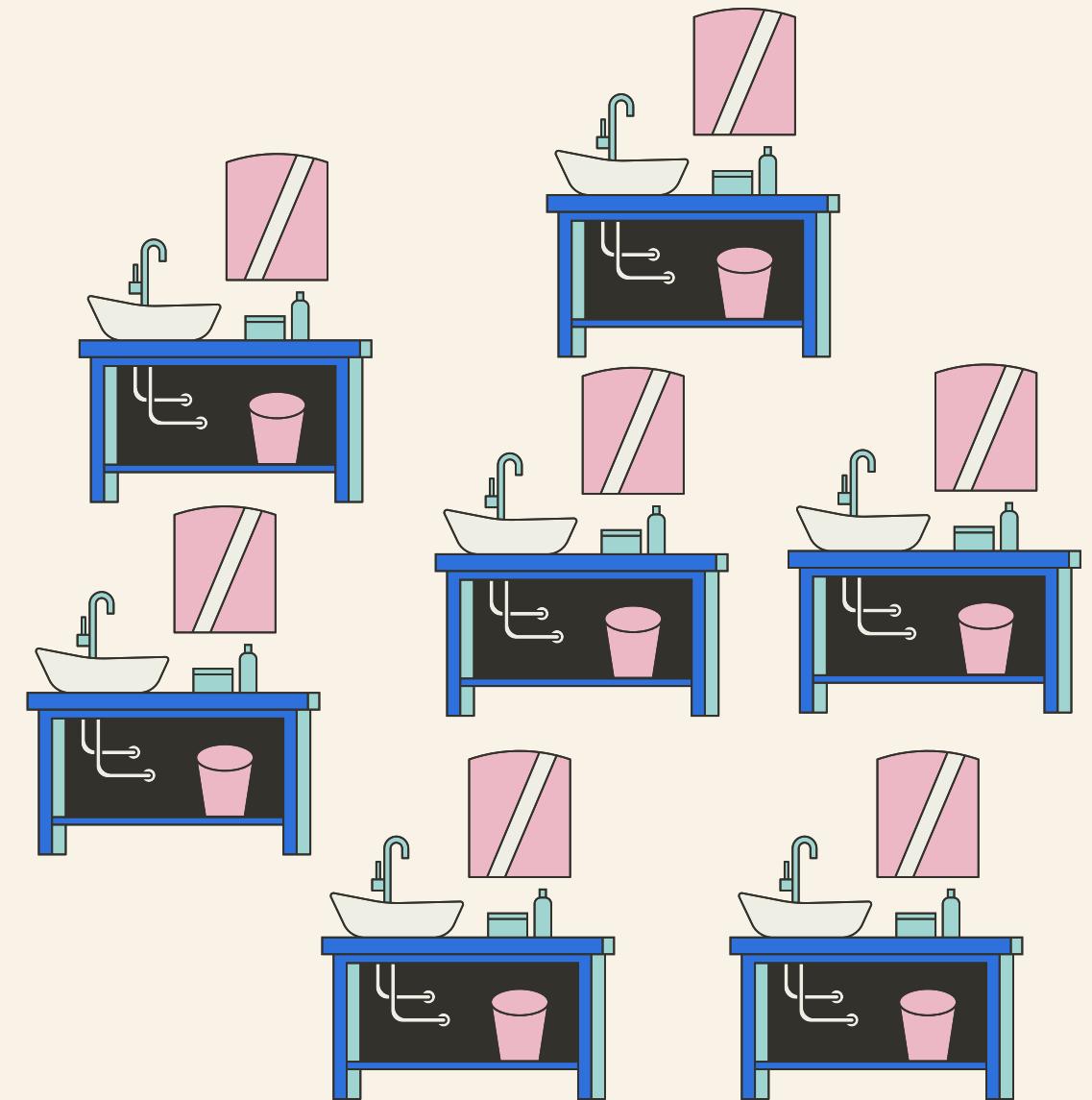
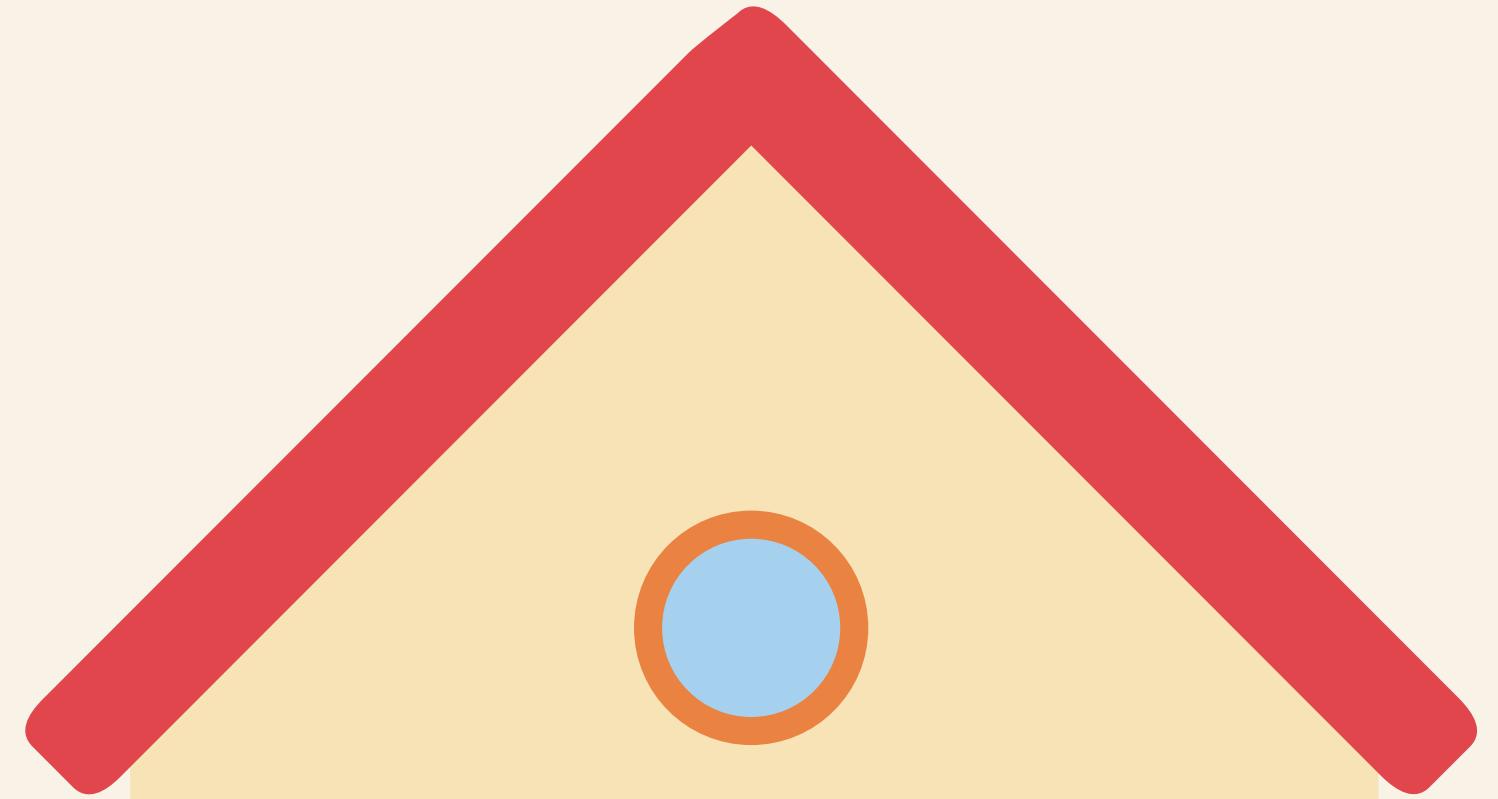
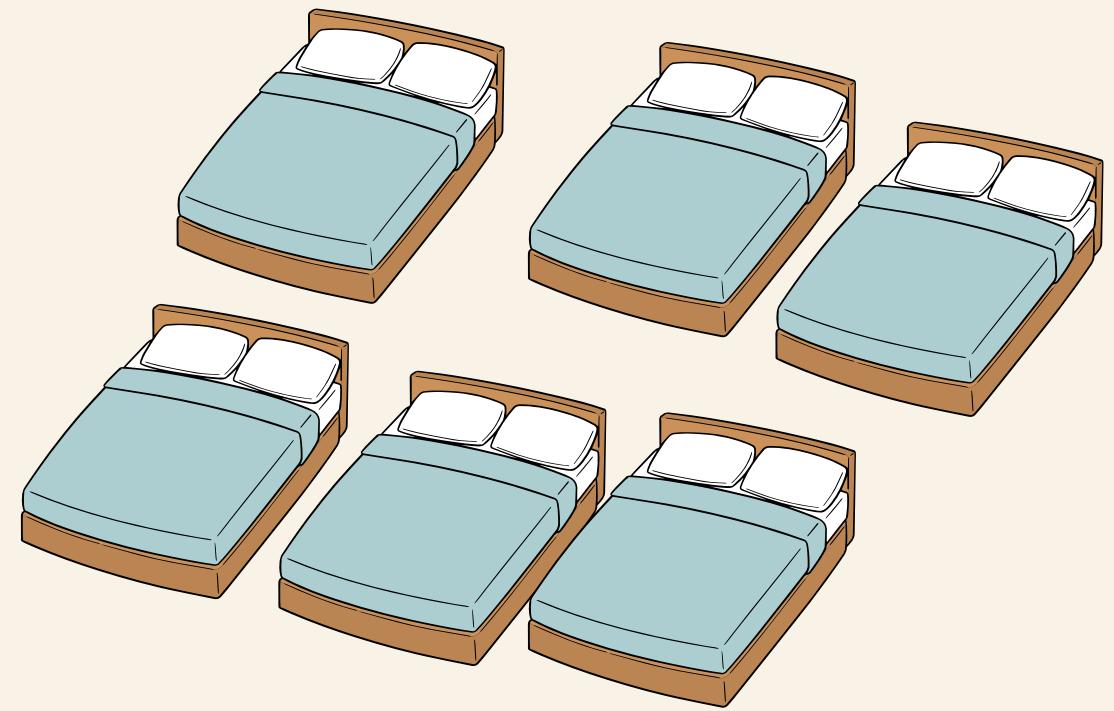




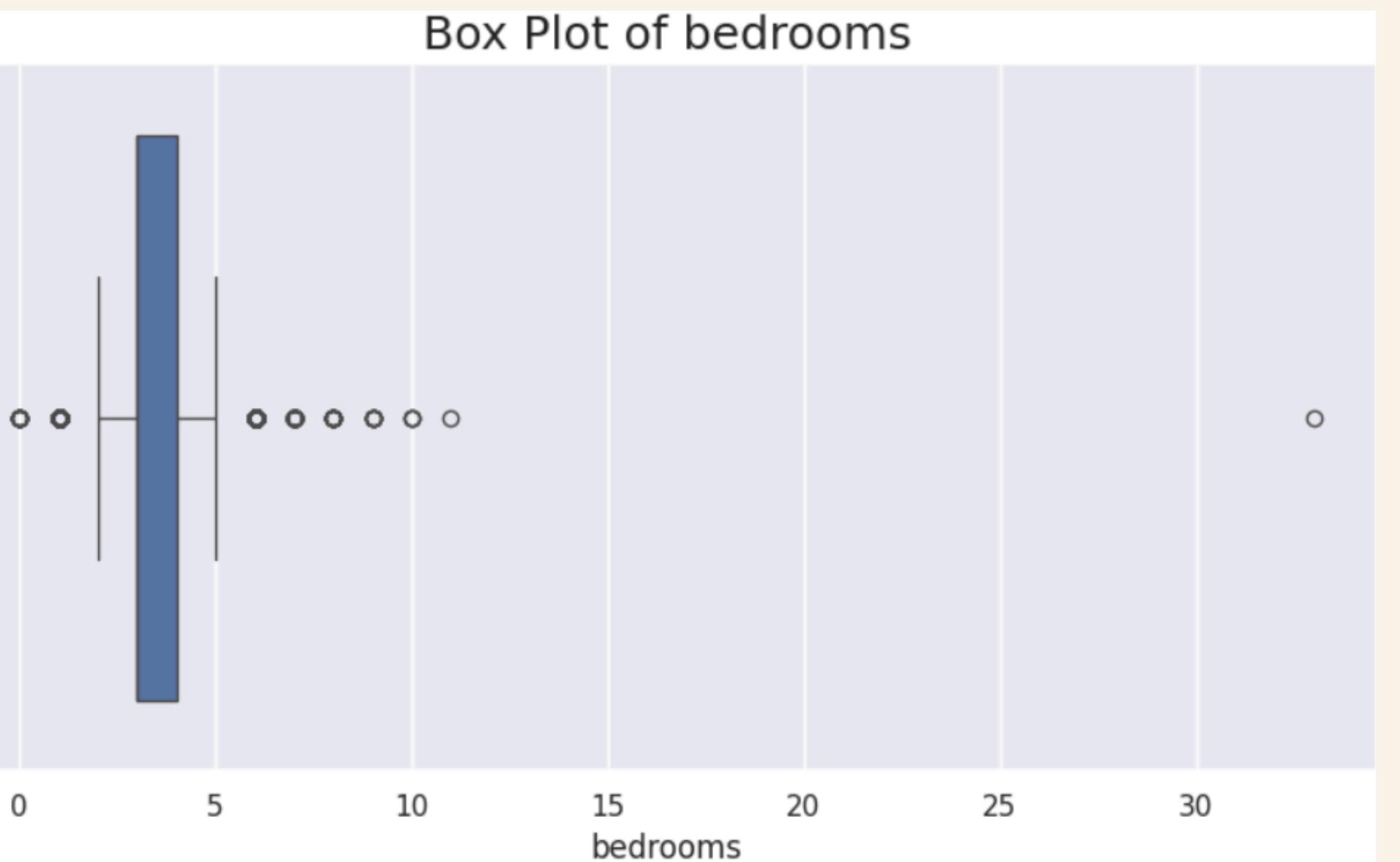
# Remove outliers

- Bigger houses are cheaper
- Depends on the owner





# Before



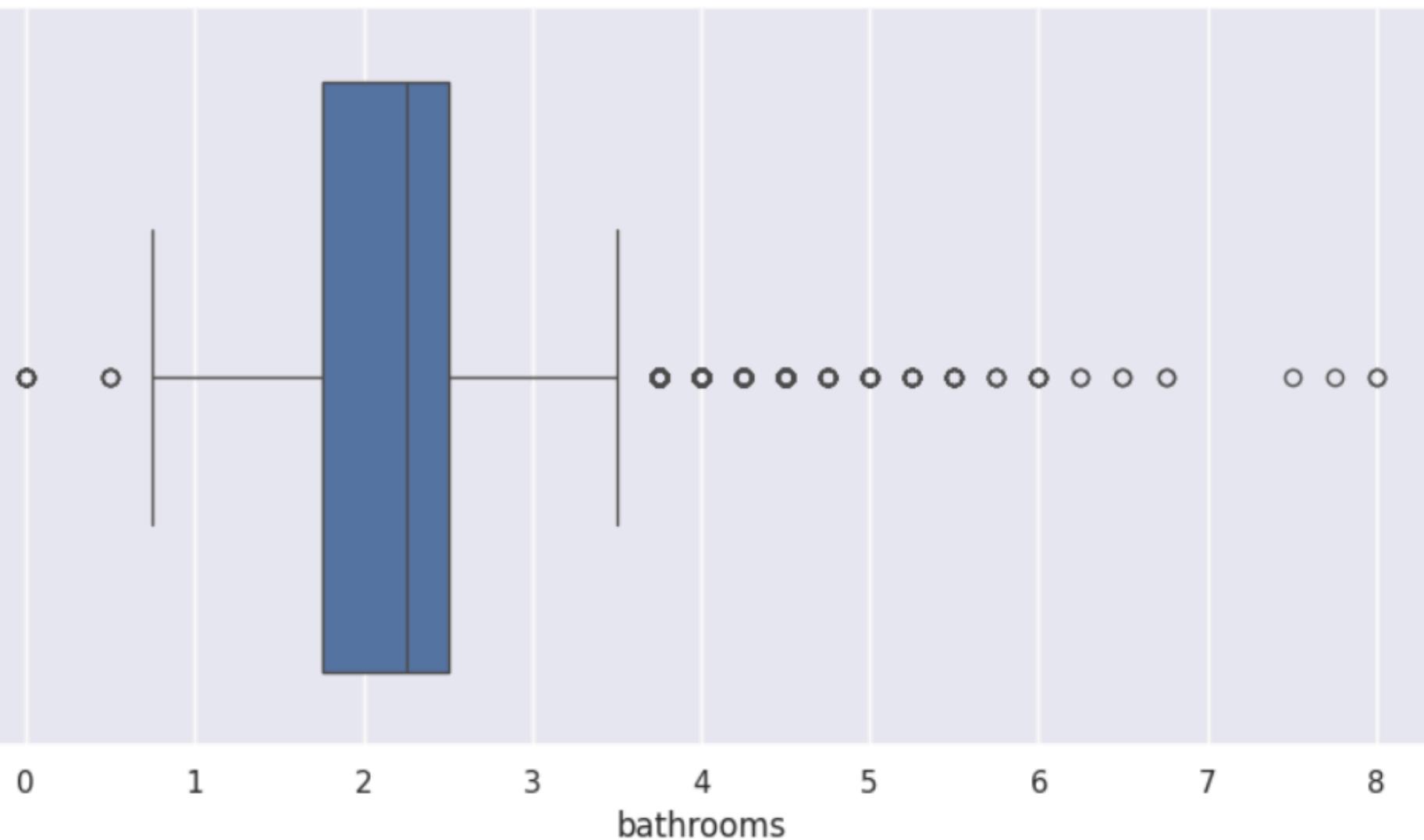
# AFTER



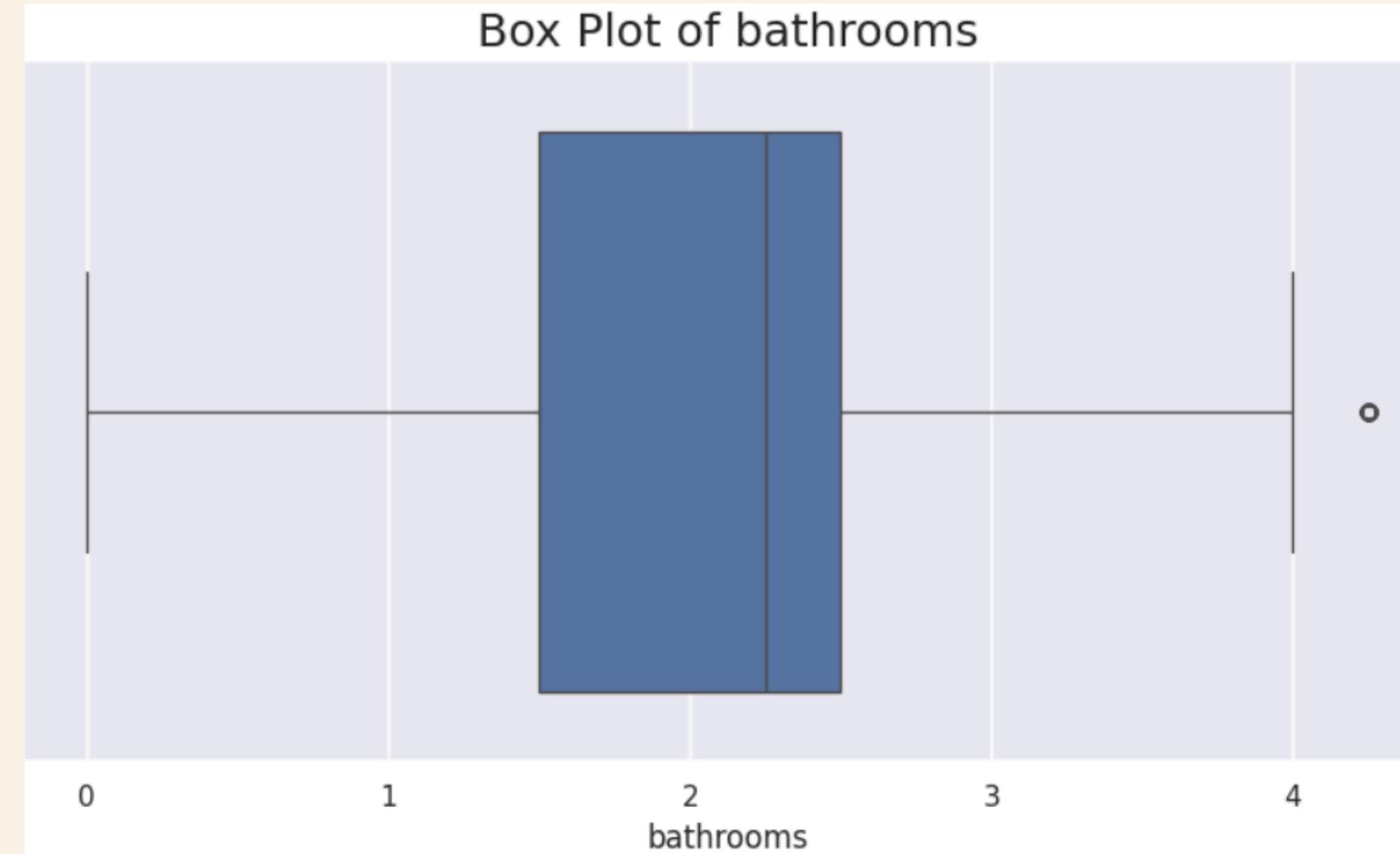
# Before

# AFTER

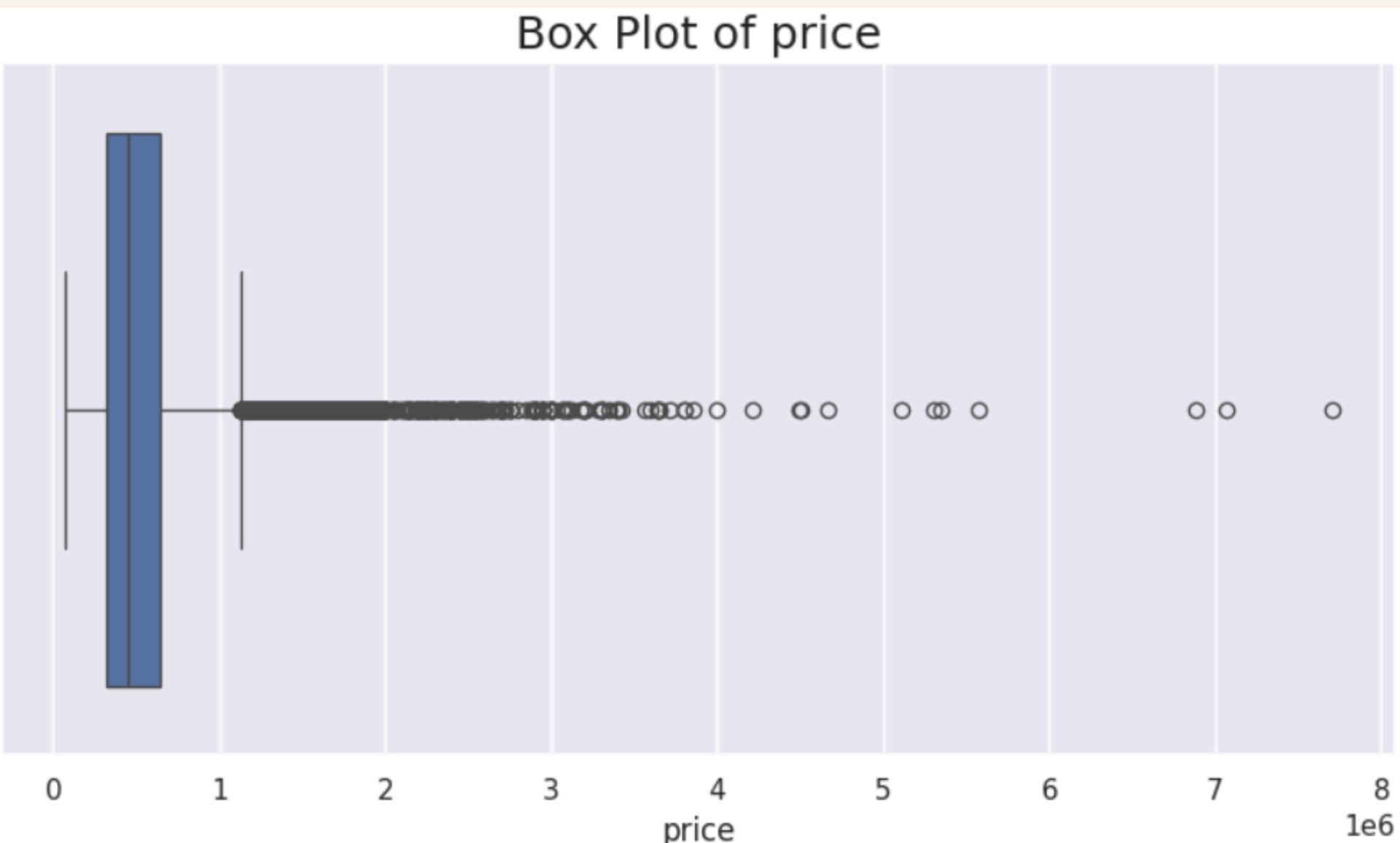
Box Plot of bathrooms



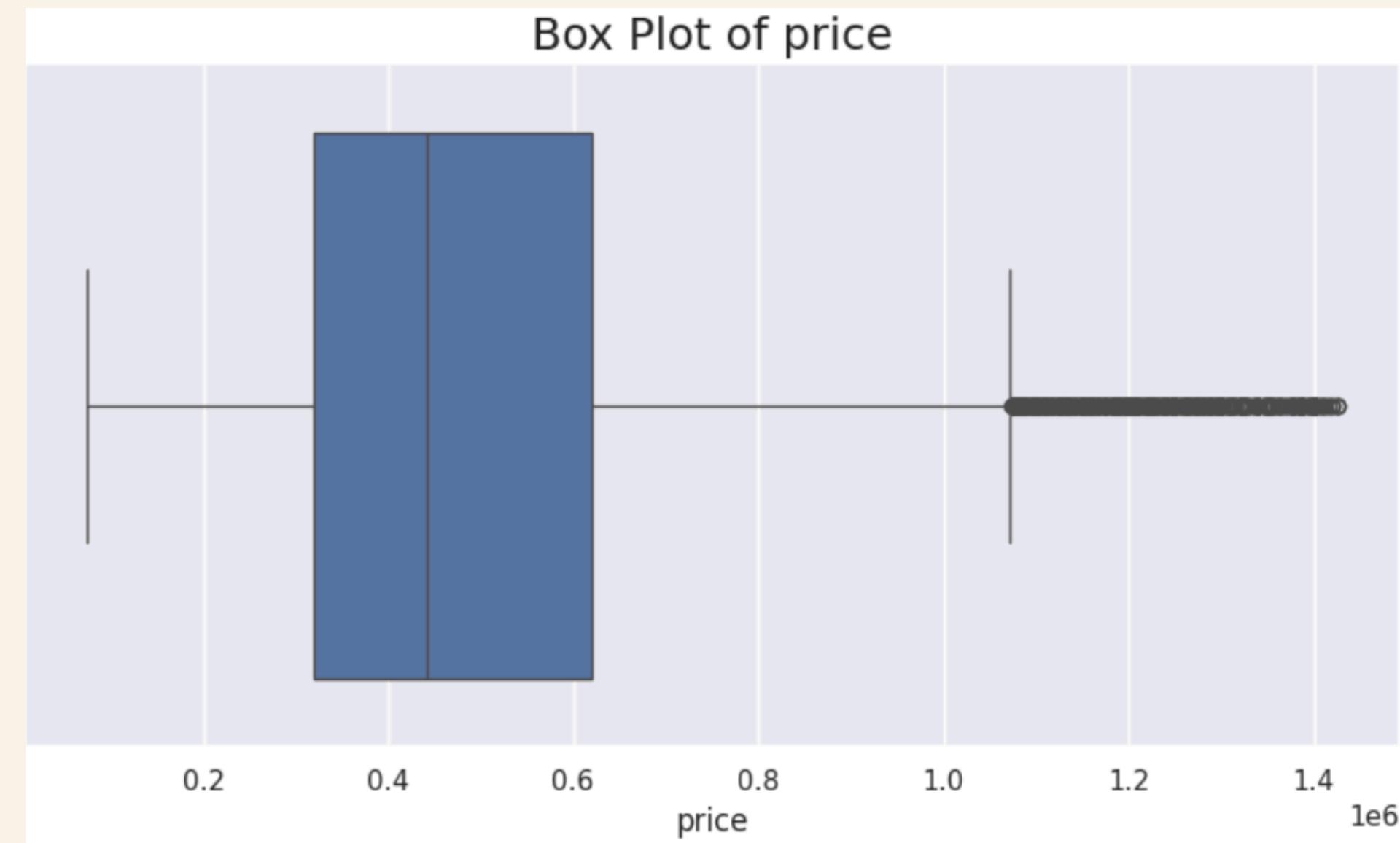
Box Plot of bathrooms



# Before



# AFTER

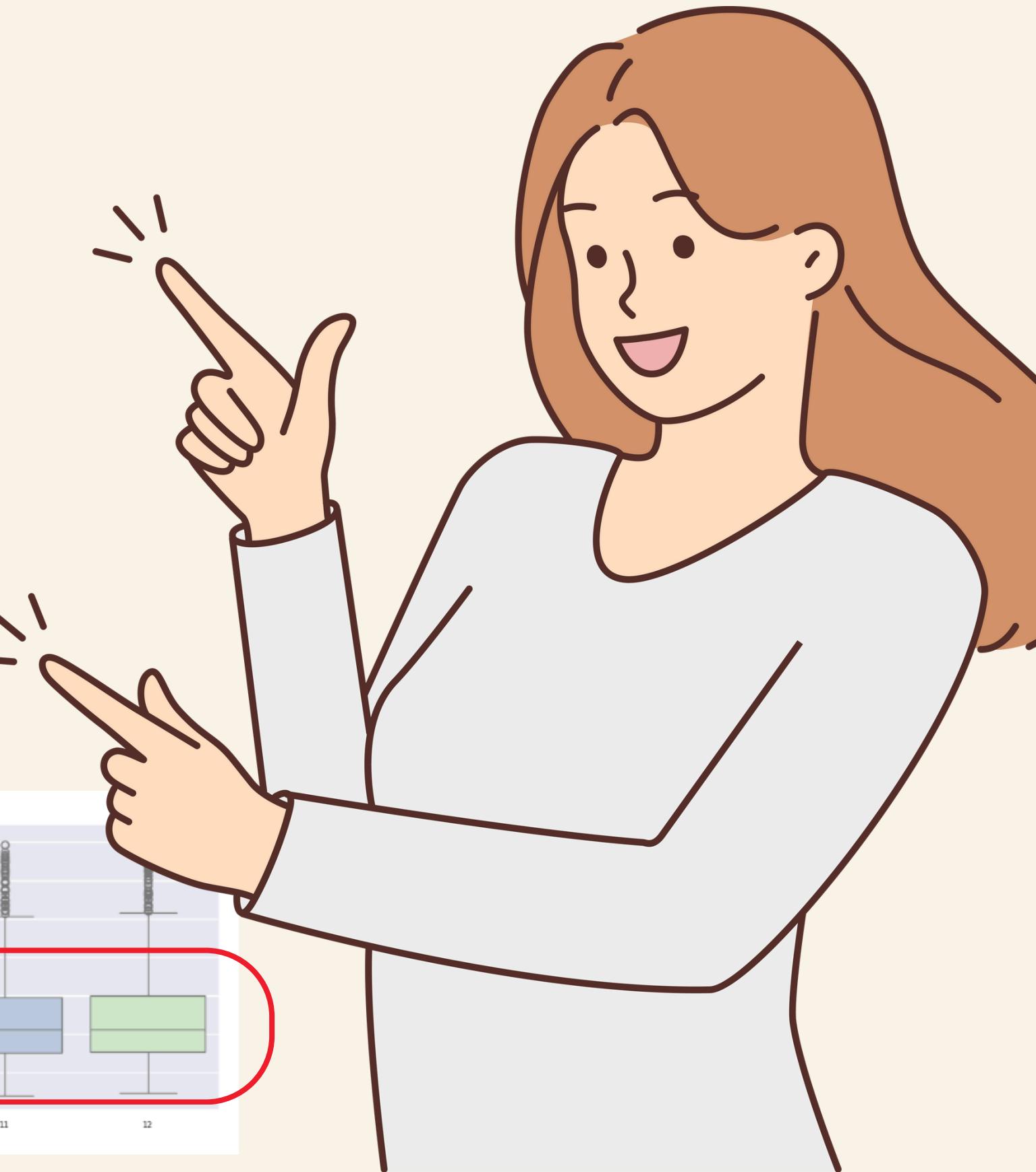
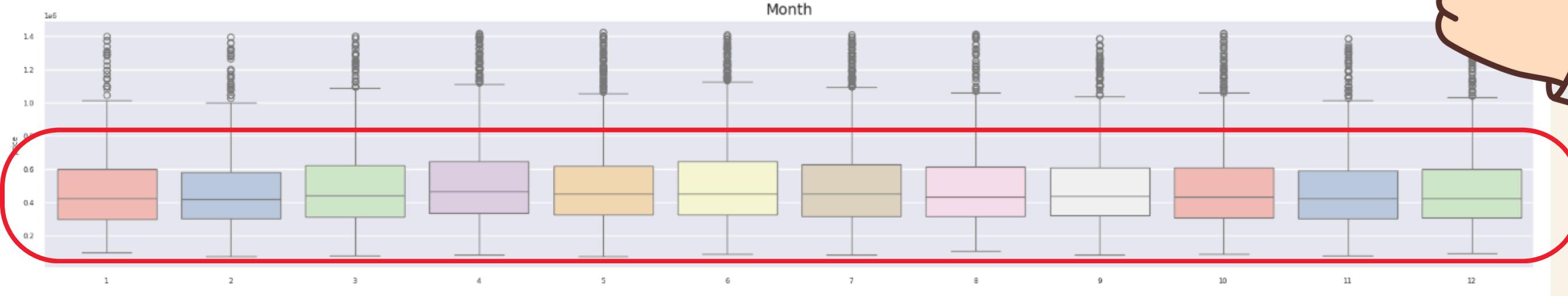
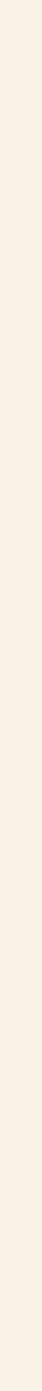
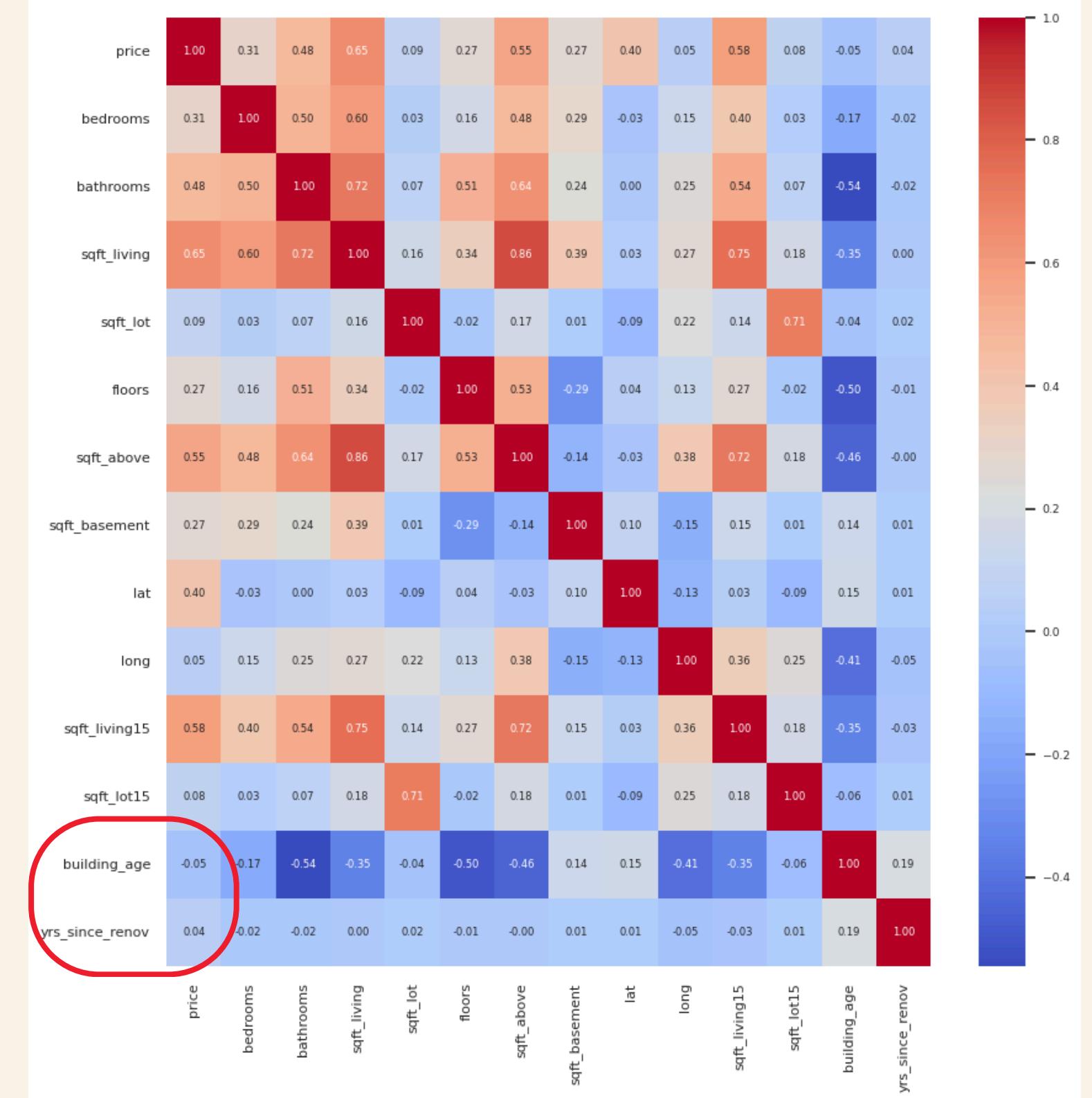


# Feature Engineering

**Age of property** = Year from “date” - “yr\_built”

**No. of years since renovated** = Year from “date” - “yr\_renovated”

**Month sold** = Month from “date”



# Feature Engineering

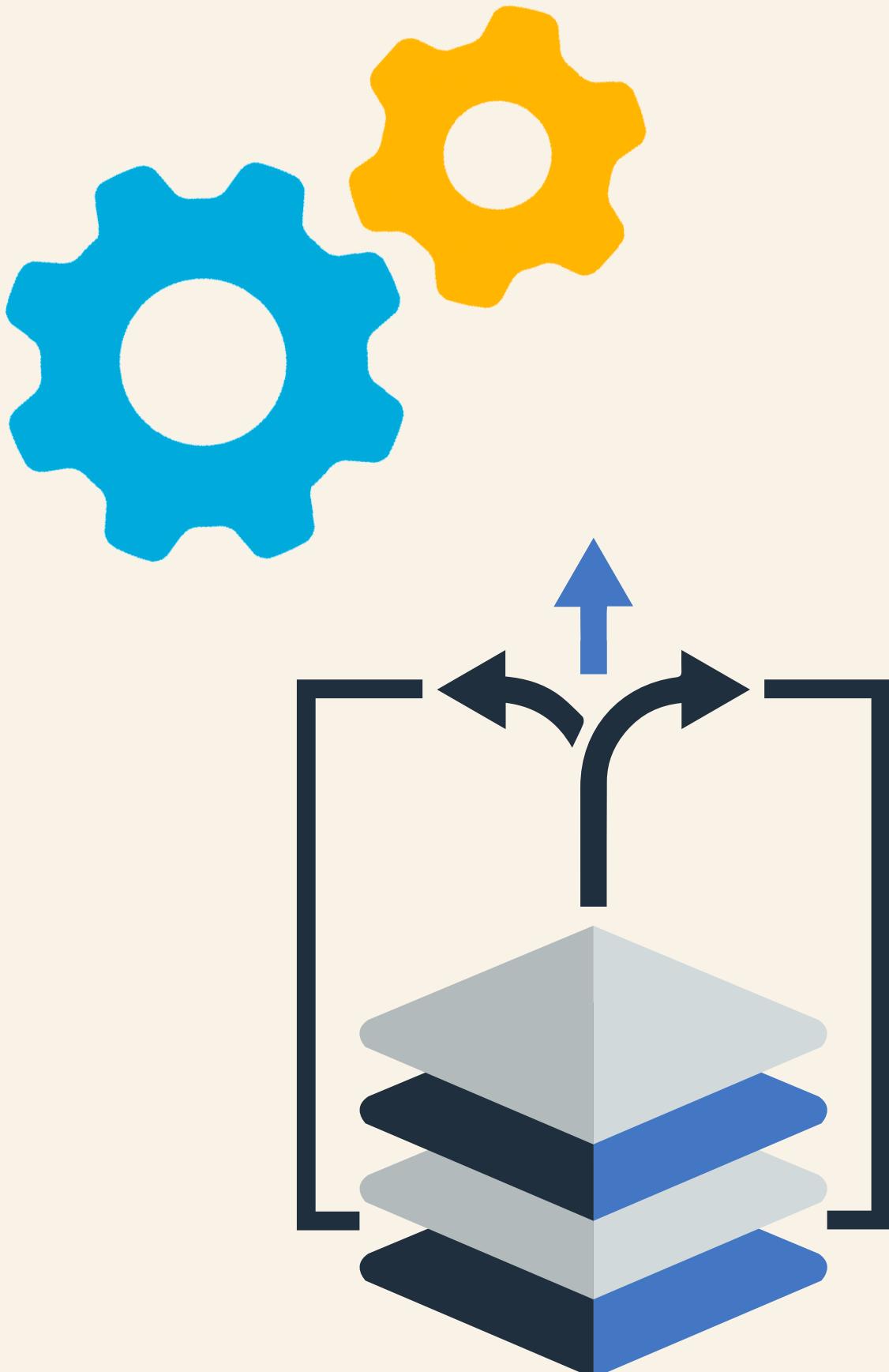
**Age of property** = Year from “date” - “yr\_built”

**No. of years since renovated** = Year from “date” - “yr\_renovated”

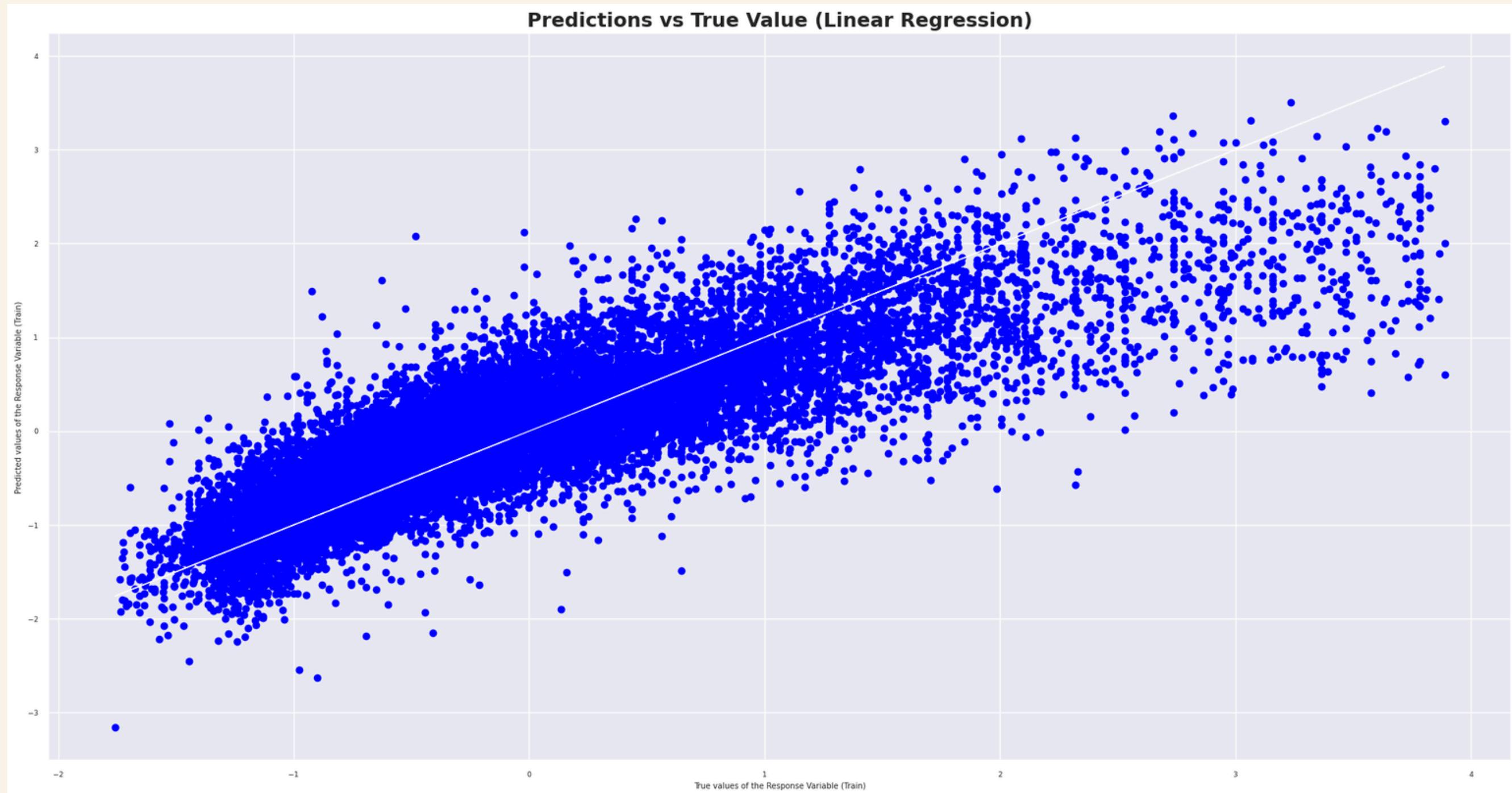
**Month sold** = Month from “date”

# Machine Learning

- Linear Regression
- Support Vector Regressor
- Decision Tree Regressor
- Random Forest Regressor.



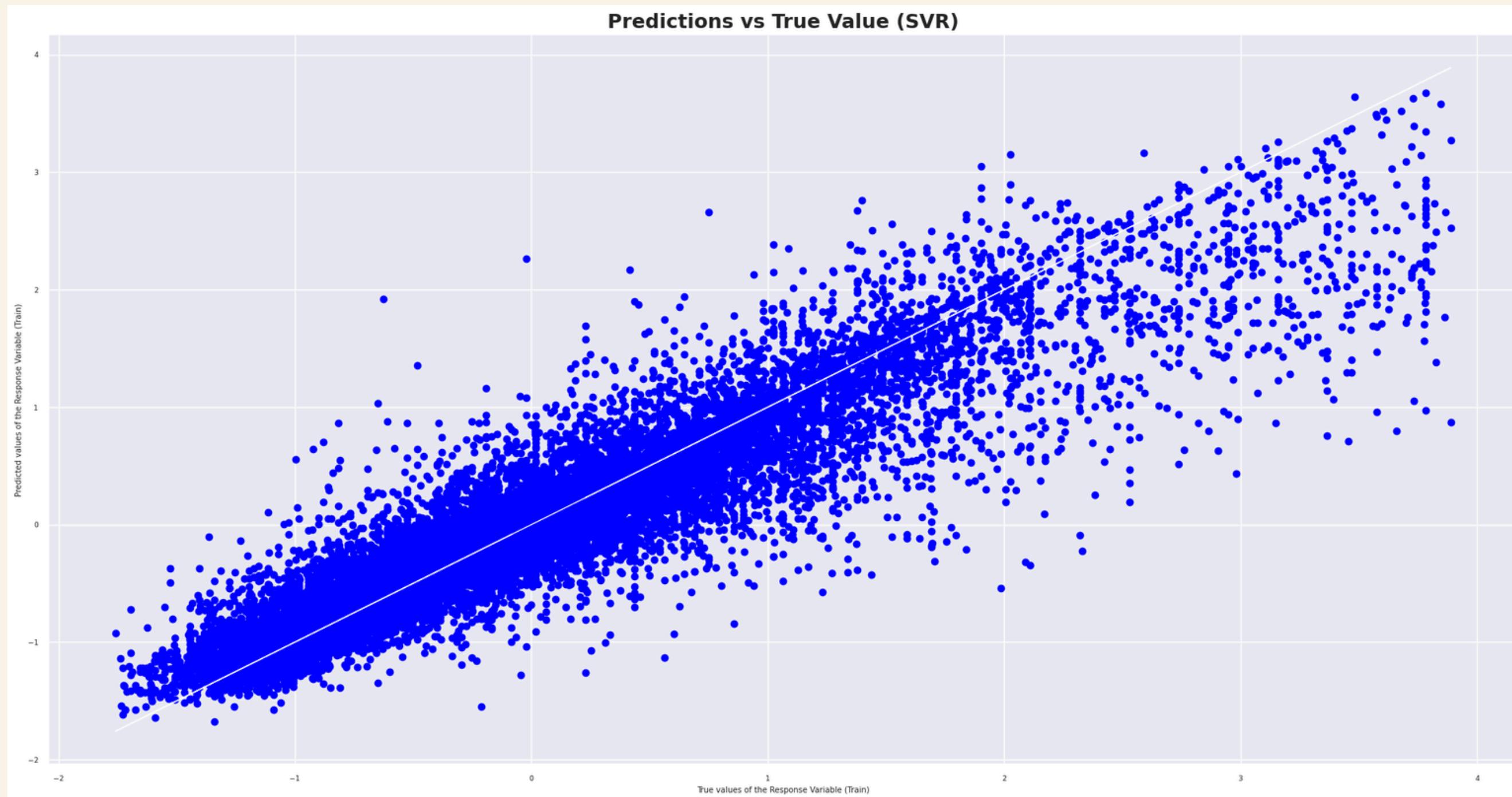
# Linear Regression



**Training accuracy  
= 0.708**

**Testing accuracy  
= 0.697**

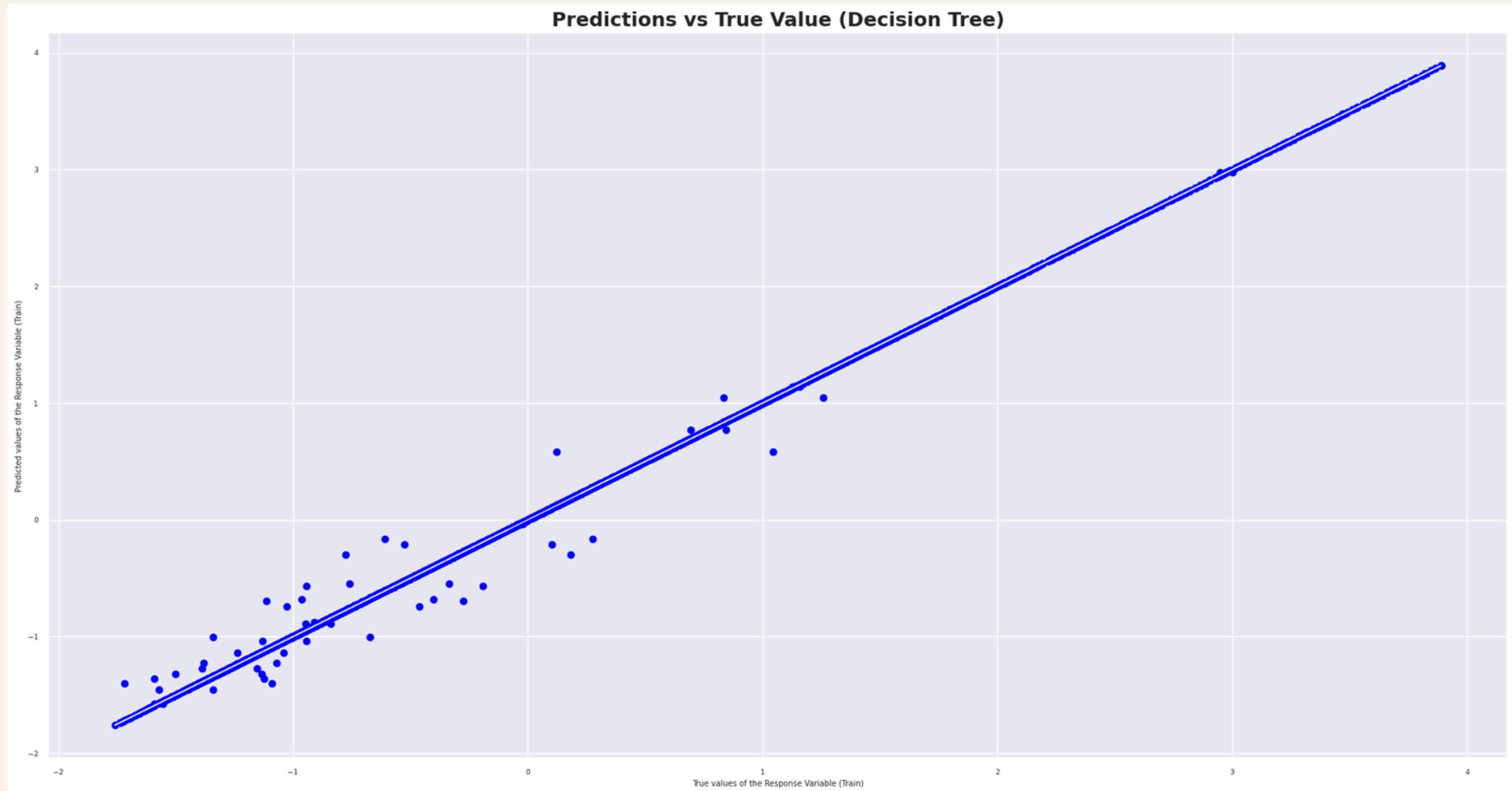
# Support Vector Regressor



**Training accuracy  
= 0.848**

**Testing accuracy  
= 0.808**

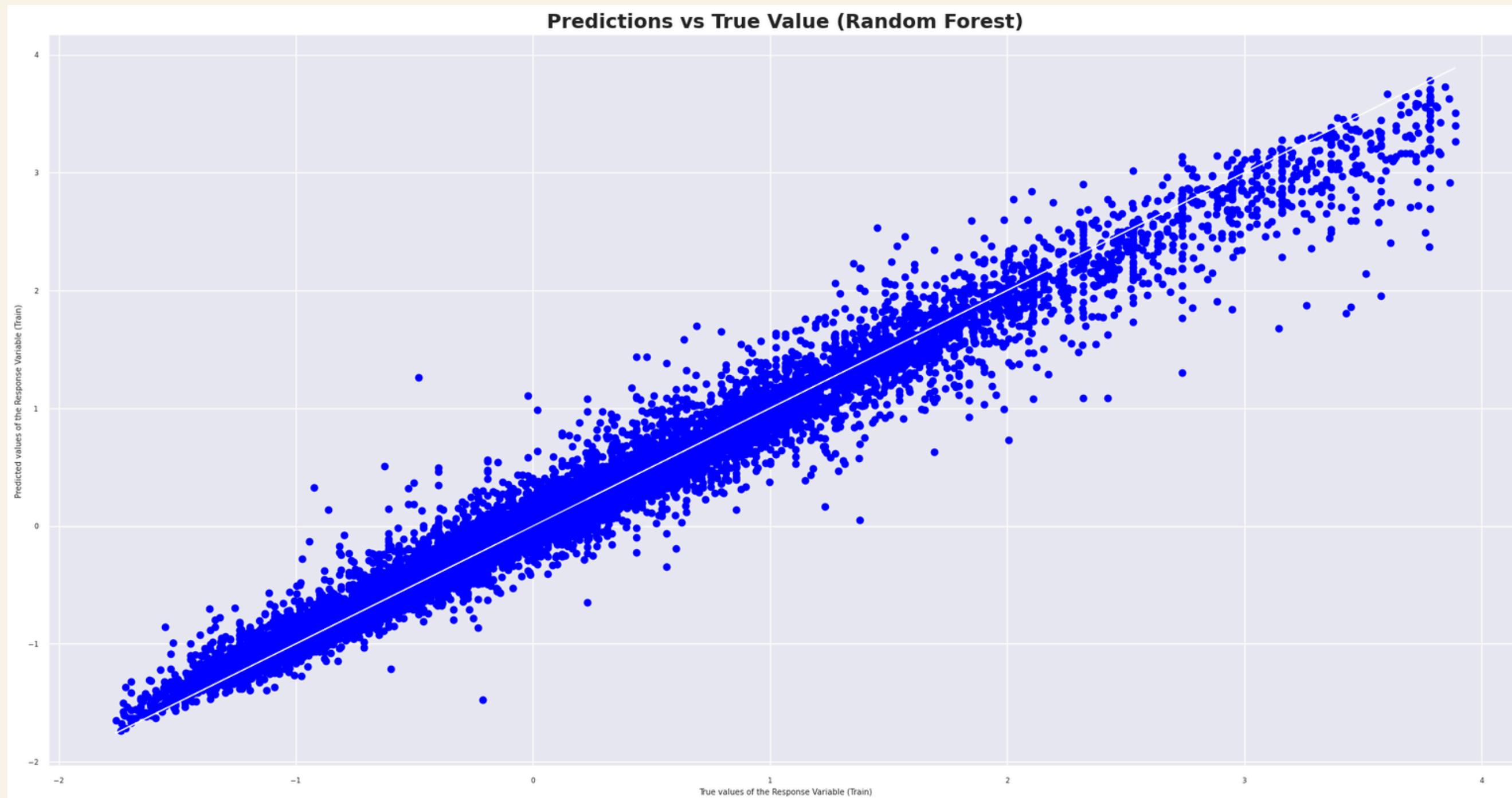
# Decision Tree Regressor



**Training accuracy**  
**= 0.999**

**Testing accuracy**  
**= 0.719**

# Random Forest Regressor



**Training accuracy**  
**= 0.972**

**Testing accuracy**  
**= 0.845**



R-squared  
Mean squared error  
Mean absolute percentage error

	Accuracy (R <sup>2</sup> score)	MSE	MAPE
Linear Regression	0.6967	0.3031	203.3716
SVR	0.8080	0.1919	149.7679
Decision Tree	0.7189	0.2810	243.3641
Random Forest	0.8453	0.1546	157.5226



	Training Error	Testing Error
<b>SVR</b>	0.1515	0.1919
<b>Random Forest</b>	0.0275	0.1546





```
Random Forest Validation R^2: 0.8507727698622572
Random Forest Validation MSE: 0.15209237961780145
Random Forest Testing R^2: 0.8676280784009862
Random Forest Testing MSE: 0.13410534487871575
```

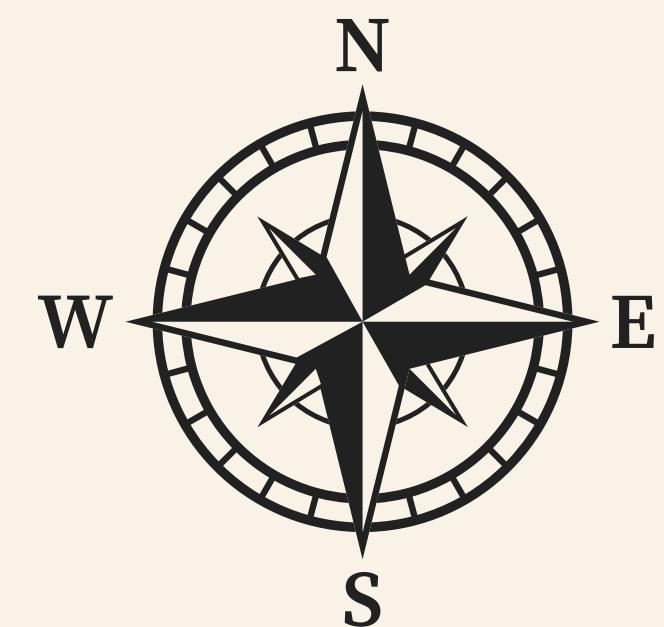
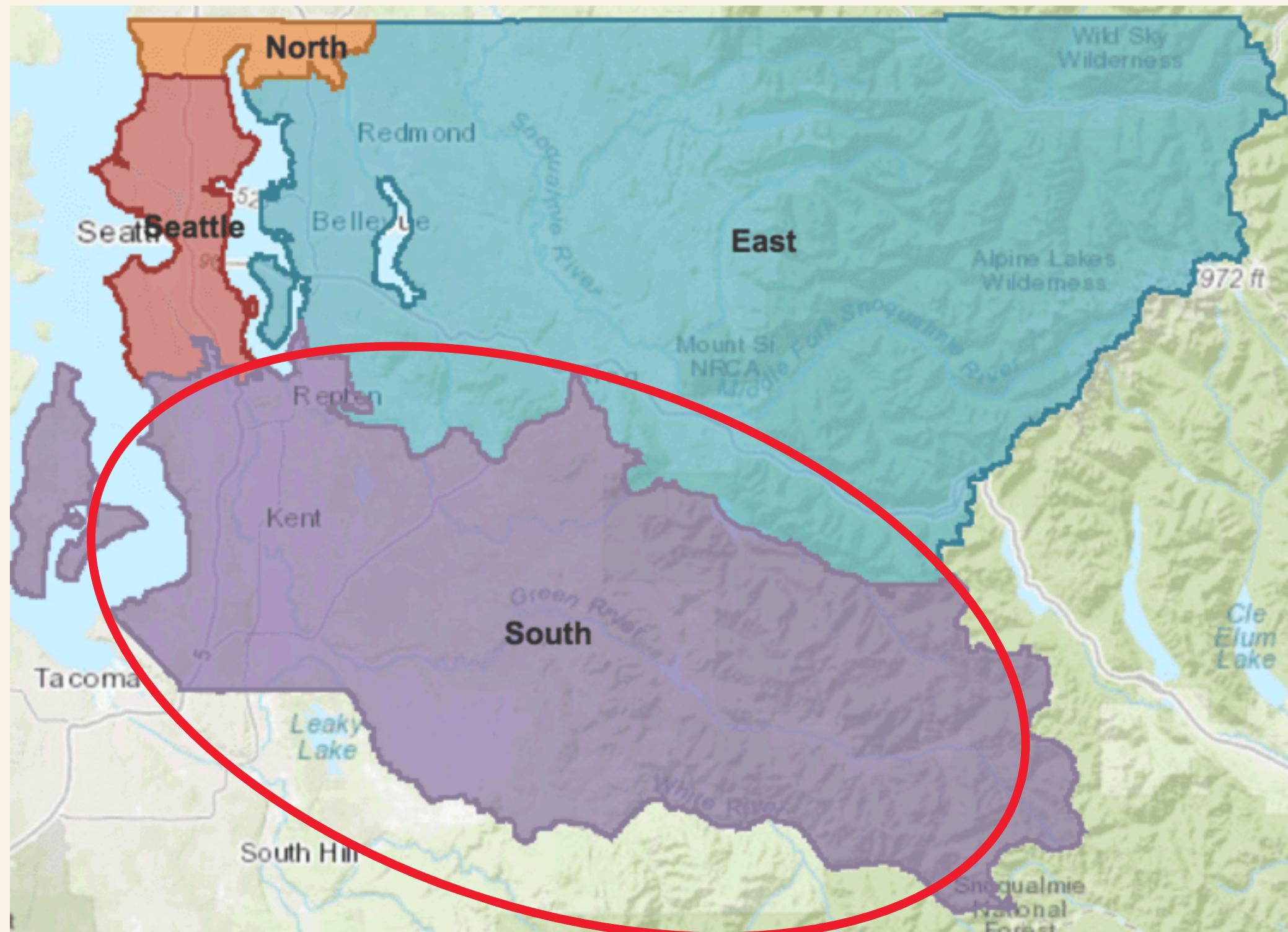


# Random Forest Regressor

# Insights and recommendations



Feature 1: grade, Score: 0.3406  
Feature 2: lat, Score: 0.3008  
Feature 3: sqft\_living, Score: 0.1613  
Feature 4: long, Score: 0.0509  
Feature 5: sqft\_living15, Score: 0.0304  
Feature 6: building\_age, Score: 0.0298  
Feature 7: sqft\_above, Score: 0.0168  
Feature 8: sqft\_lot, Score: 0.0168  
Feature 9: view, Score: 0.0155  
Feature 10: sqft\_lot15, Score: 0.0153





373818.01

(Predicted by best model)





