# Ames Housing Price Prediction
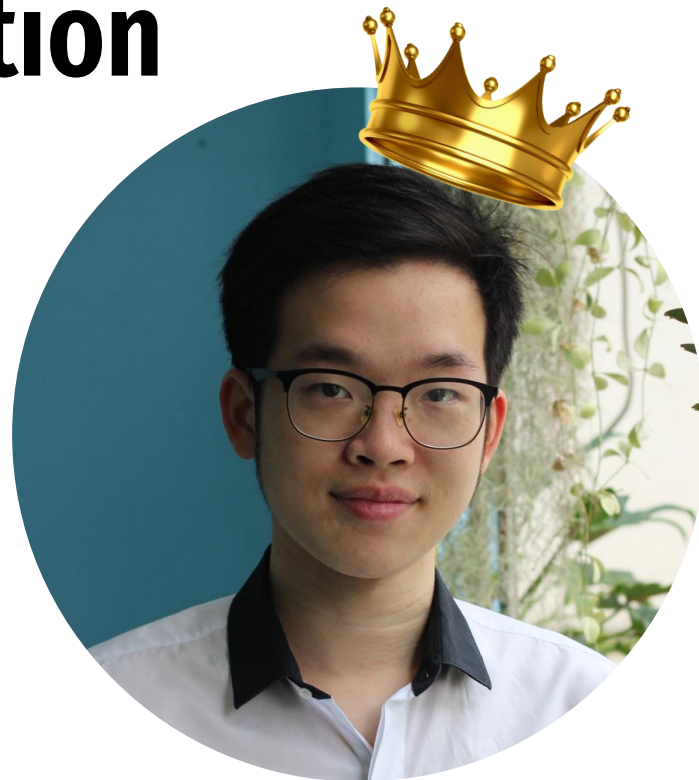
**More Power!**

**By Im Depends**

# Introduction

Chalermchon Wongsopa

Kantaphon Vareekasem

# Problem Statement

**How much?**



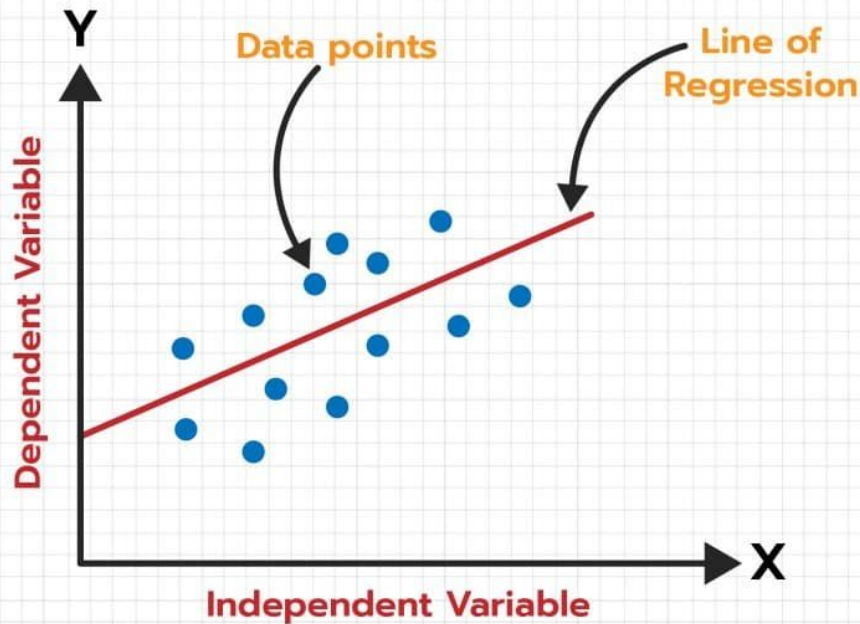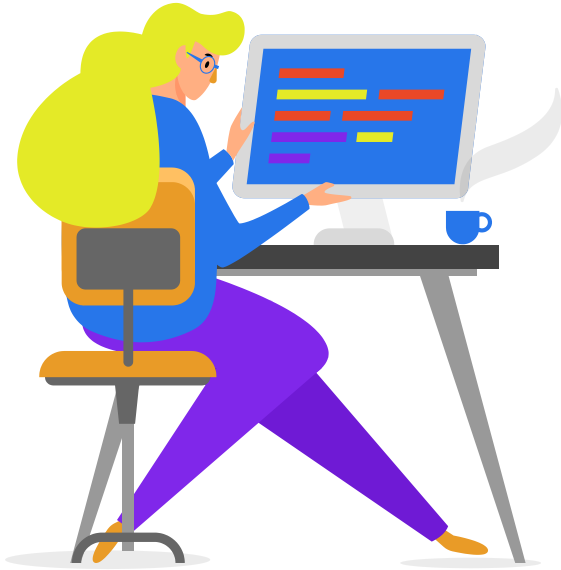| **Information** | Many Features |
| **Average Price** | Average price can only give rough estimate. |
| **Overprice** | Avoid overpriced house |
| **Underprice** | Opportunities to make money |

# Solution

## Linear Regression



### Best Model

15 Numerical + 20 Dummies

**RMSE**: **19,828** USD

**R Squared**: **93%**

# Models Improvement Process



**Baseline Model**

01

Top 10 Features

**RMSE: 36,082**

**Core Model**

02

Top 15 Features

**RMSE: 34,478**

**Outliers**

03

Cleaning

**RMSE: 29,417**

**Dummy Variables**

04

Categorical Variables

**RMSE: 23,351**
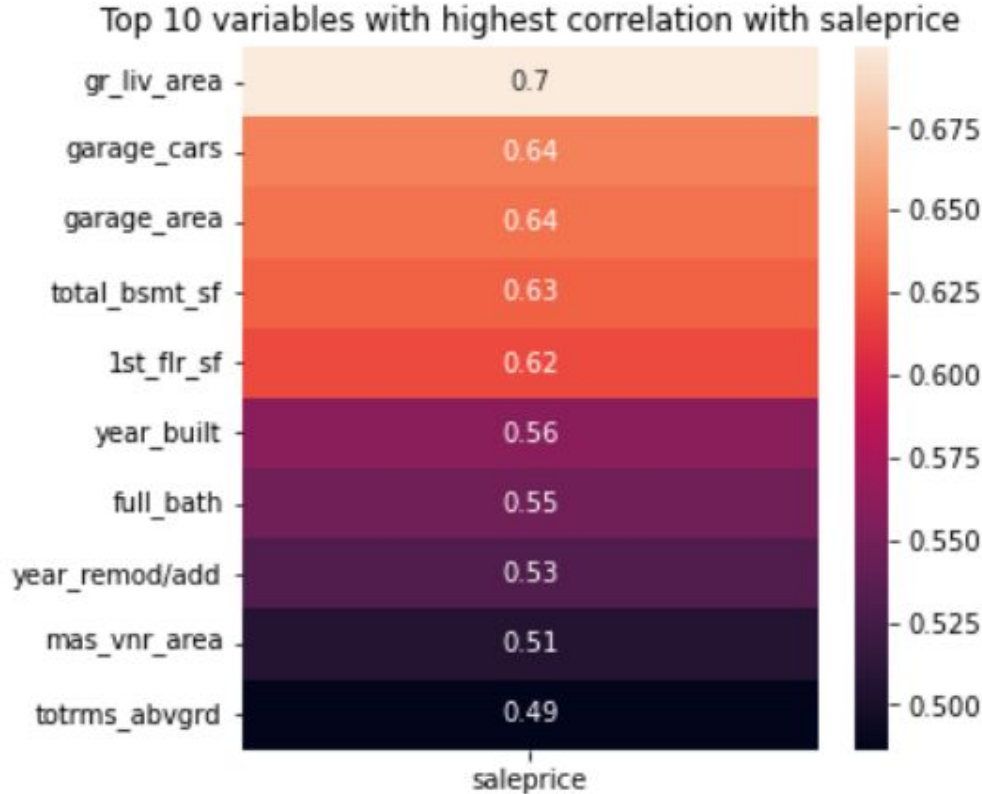
**Log Sale Price**

05

Log Y

**RMSE: 21,023**

**Others**

06
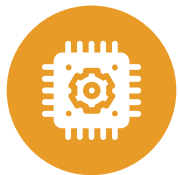
Interaction term,
Log X, Train More

**RMSE:19,828**

# Baseline Model



Top 10 variables with highest correlation with saleprice

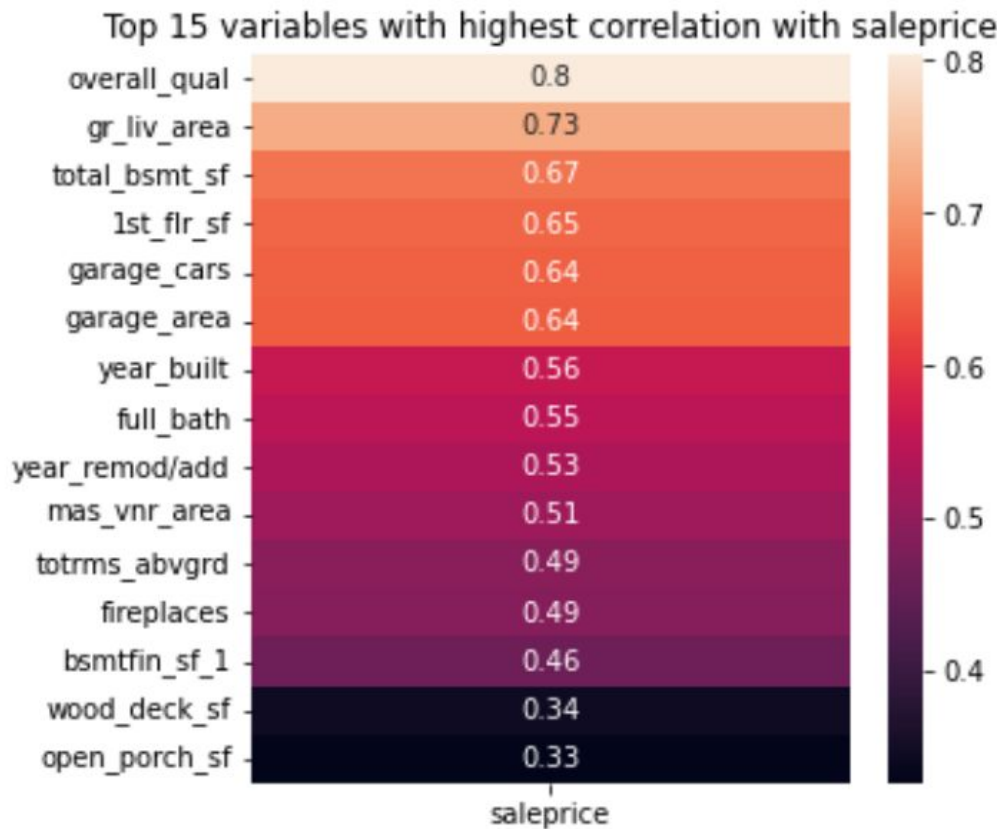| | saleprice |
|---|---|
| gr_liv_area | 0.7 |
| garage_cars | 0.64 |
| garage_area | 0.64 |
| total_bsmt_sf | 0.63 |
| 1st_flr_sf | 0.62 |
| year_built | 0.56 |
| full_bath | 0.55 |
| year_remod/add | 0.53 |
| mas_vnr_area | 0.51 |
| totrms_abvgrd | 0.49 |

**Linear Regression**
- Top 10 variables
- Fill Missing Values with 0
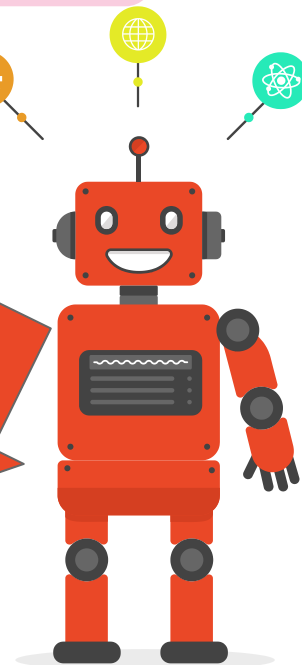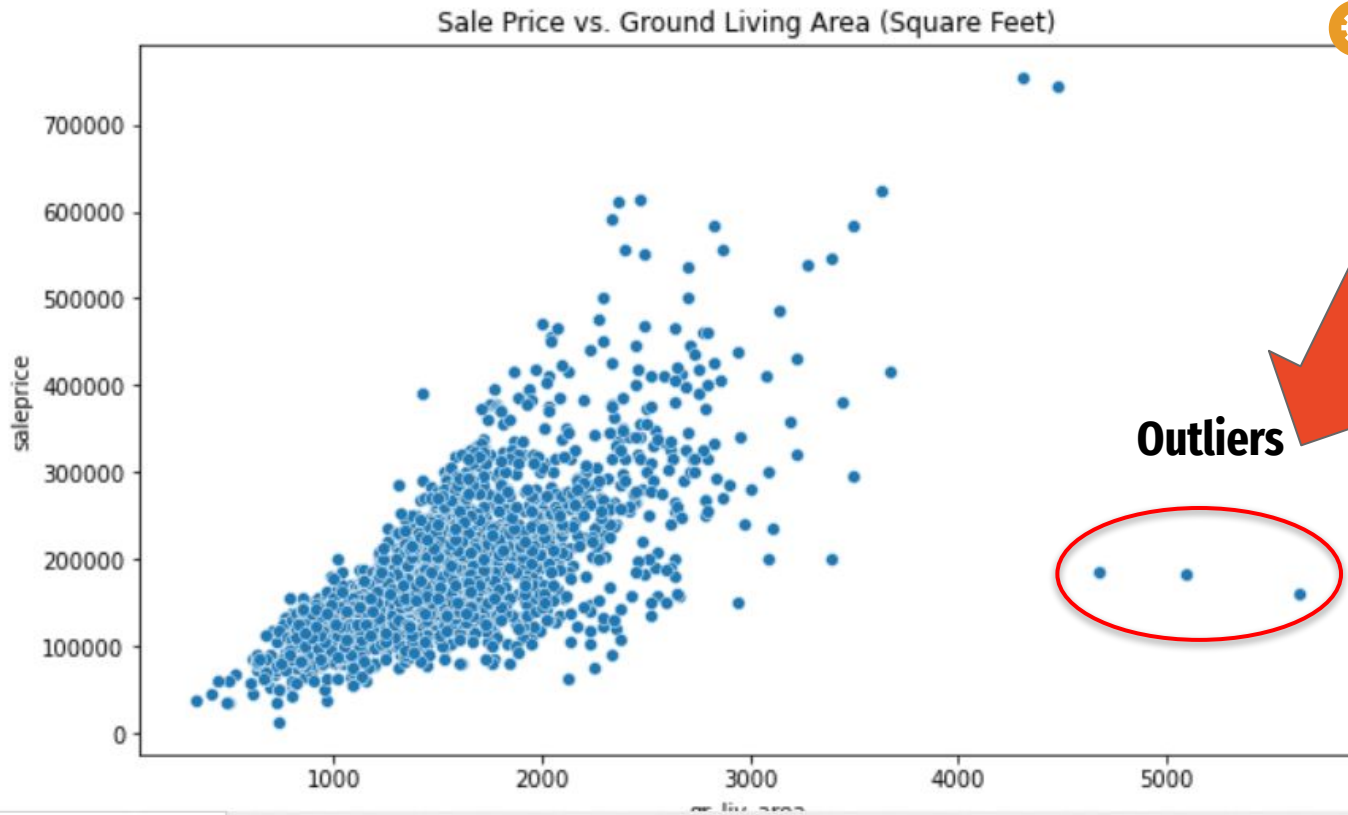- Target Variable: Sale Price

**RMSE: 36,082**

# Core Model



Top 15 variables with highest correlation with saleprice

| Variable | saleprice |
|---|---|
| overall_qual | 0.8 |
| gr_liv_area | 0.73 |
| total_bsmt_sf | 0.67 |
| 1st_flr_sf | 0.65 |
| garage_cars | 0.64 |
| garage_area | 0.64 |
| year_built | 0.56 |
| full_bath | 0.55 |
| year_remod/add | 0.53 |
| mas_vnr_area | 0.51 |
| totrms_abvgrd | 0.49 |
| fireplaces | 0.49 |
| bsmtfin_sf_1 | 0.46 |
| wood_deck_sf | 0.34 |
| open_porch_sf | 0.33 |

**Linear Regression**
- Top 15 variables
- Adjust number of top correlated variables using RMSE

**RMSE: 34,478**

# Outliers

RMSE: 29,417

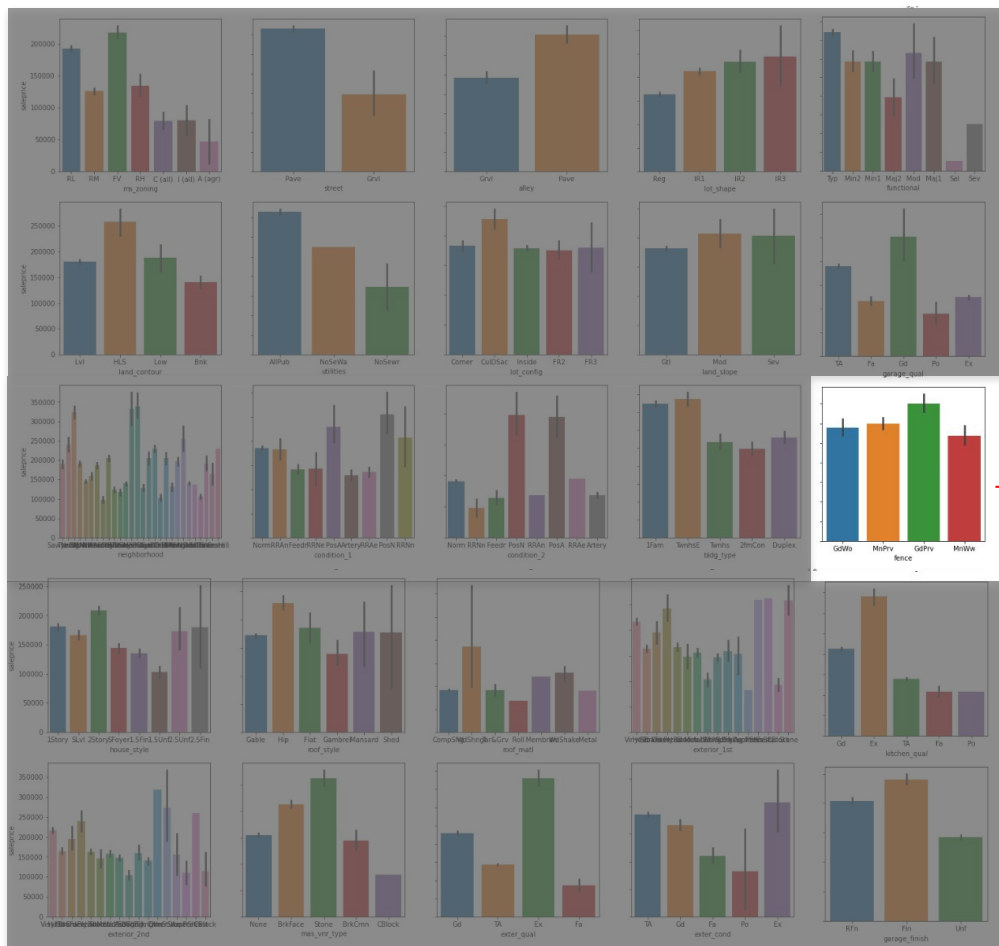Sale Price vs. Ground Living Area (Square Feet)

Outliers

# Categorical Variables



- 43 Variables are available for the model

- We distinguished 'Signal' from 'Noise' using 2 Indicators
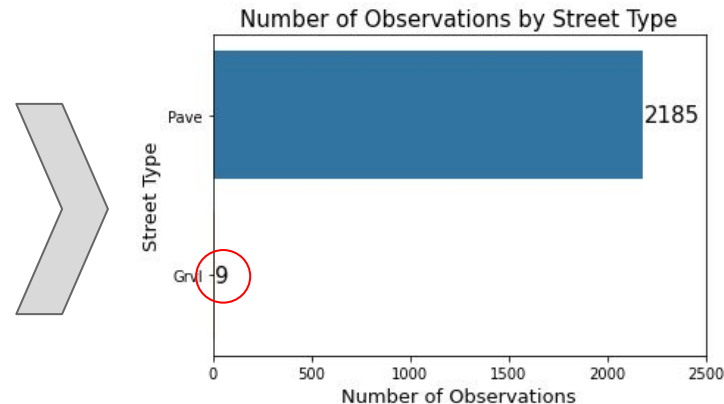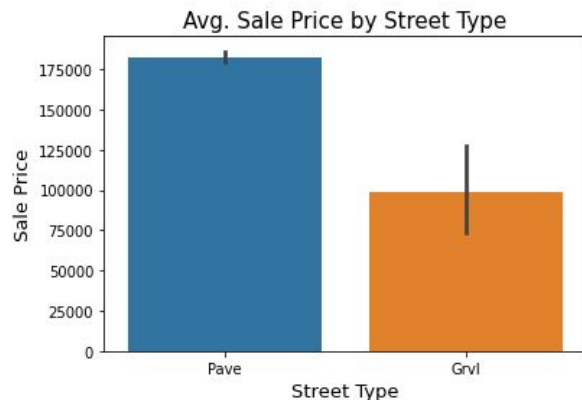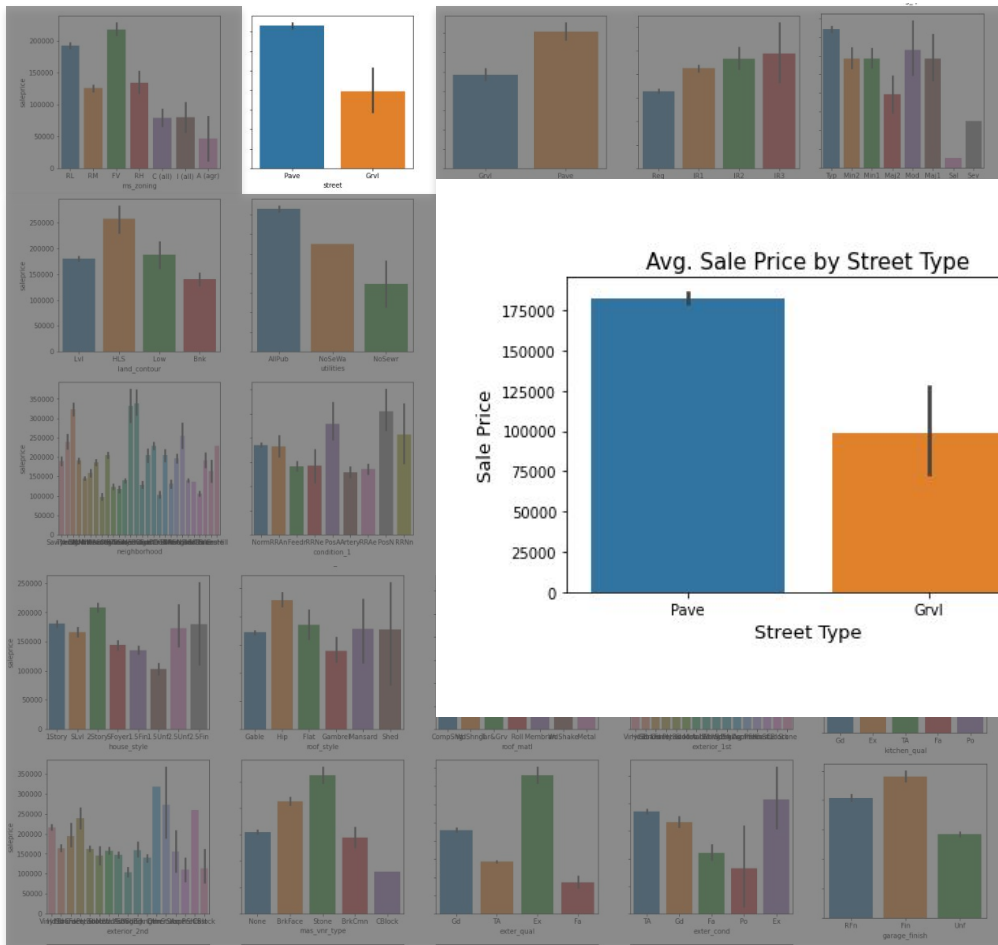  - Average sale price
  - Number of observations

**RMSE: 23,351**

# Excluded Variables - Type I



**Indistinguishable "Sale Price"** for each "Fence Quality"

# Excluded Variables - Type II



Avg. Sale Price by Street Type

Number of Observations by Street Type

**Insufficient data** for **gravel street type** to train the model

# Selected Categorical Variables

## Group I

**01**

**Applied "get_dummies" function**

1. Locations within Ames city
2. Type of dwelling
3. Exterior quality
4. Condition of sale
5. Fireplace quality
6. Flatness of the property
7. Home functionality
8. General shape of property
9. Paved driveway
10. Central air conditioning

## Group II

**02**

**Grouped to binary form**

1. Kitchen quality
2. Height of the basement
3. Type of sale
4. Exterior covering on house
5. Heating quality
6. Zoning classification
7. Proximity to various conditions
8. Garage condition
9. Garage quality
10. Miscellaneous feature

# 👍 1st Most Impact - Neighborhood



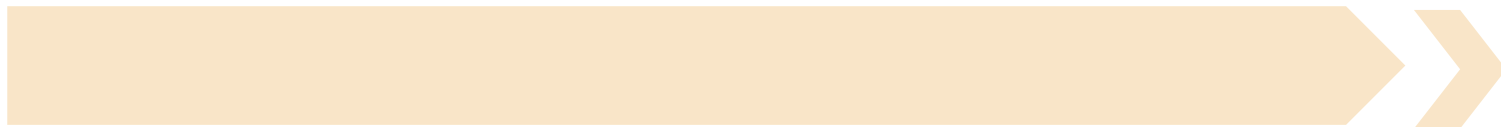Sale Price by Neighborhood

Around **1,400 USD of RMSE was decreased** after including the variable to the model

# 👍 2nd Most Impact – Type of Dwelling



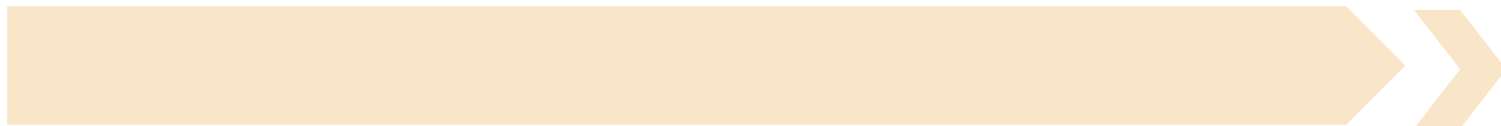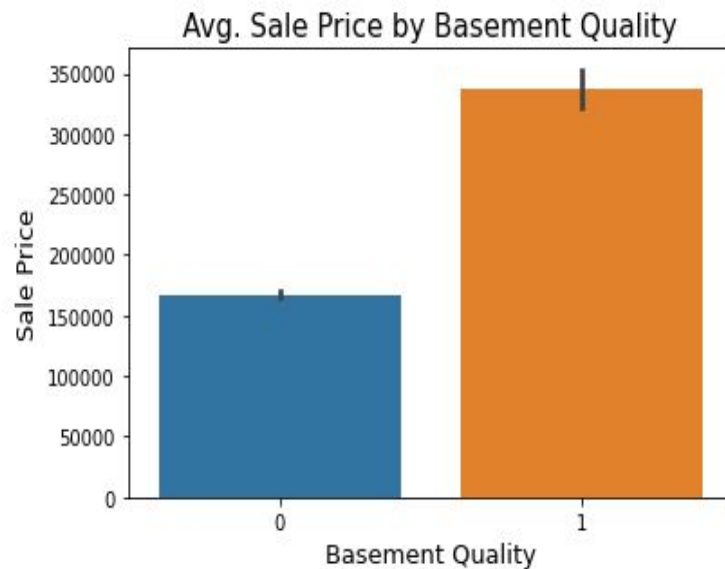Sale Price by Type of Dwelling

# 👍 1st Most Impact - Kitchen Quality
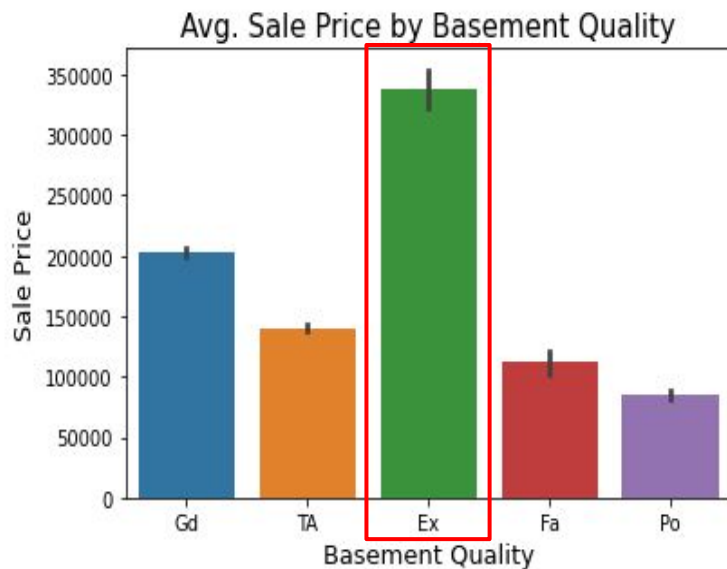


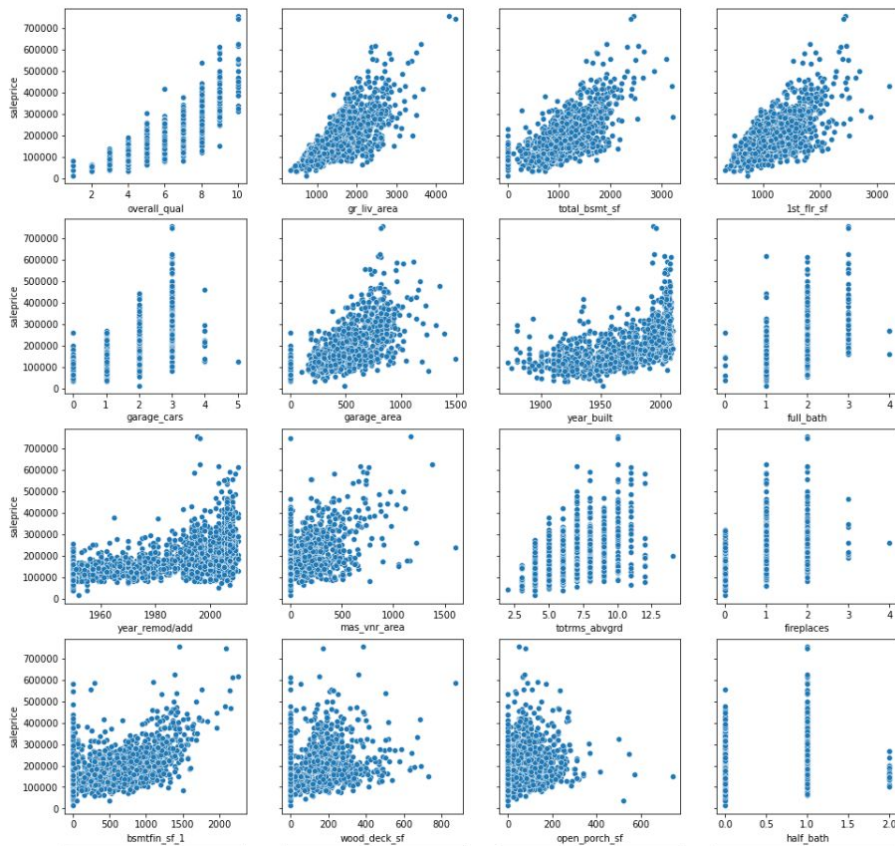Avg. Sale Price by Kitchen Quality



Avg. Sale Price by Kitchen Quality
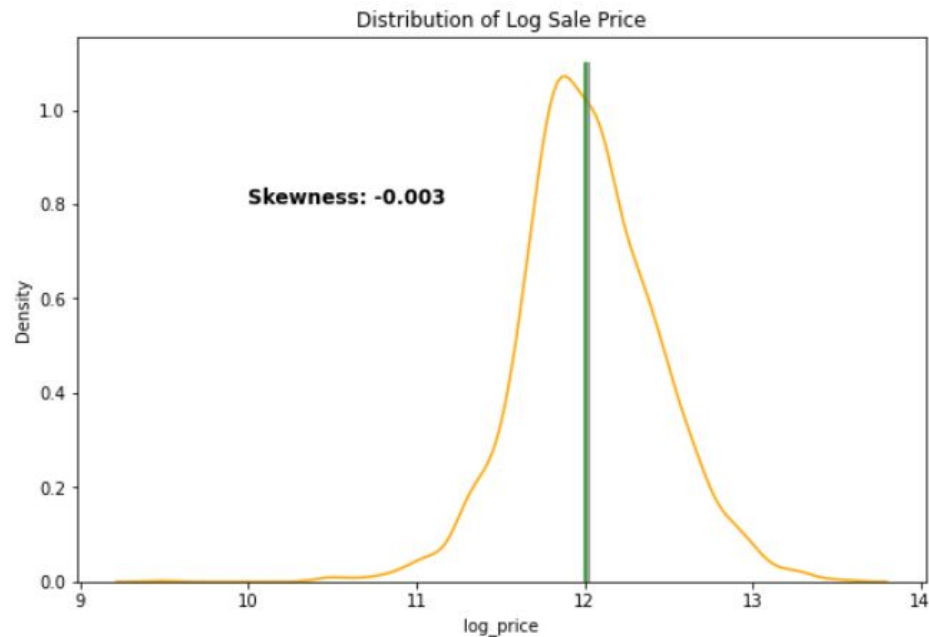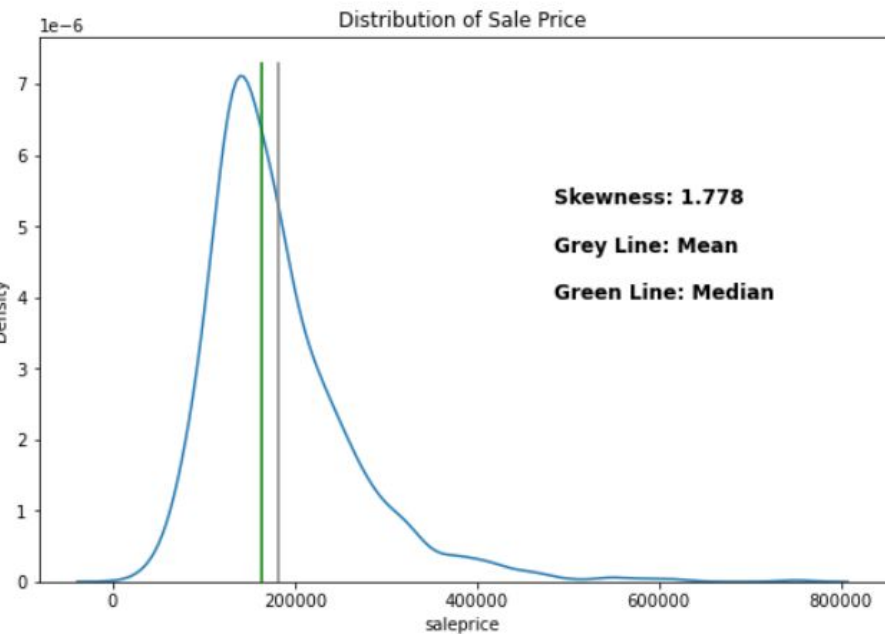
# 2nd Most Impact - Basement Quality

# Scatter Plot



**Scatter Plot**
- Check LINE ASSUMPTION
- Linearity
- If not linear, we can transform
- Drop variables: Open-Porch, Total Room
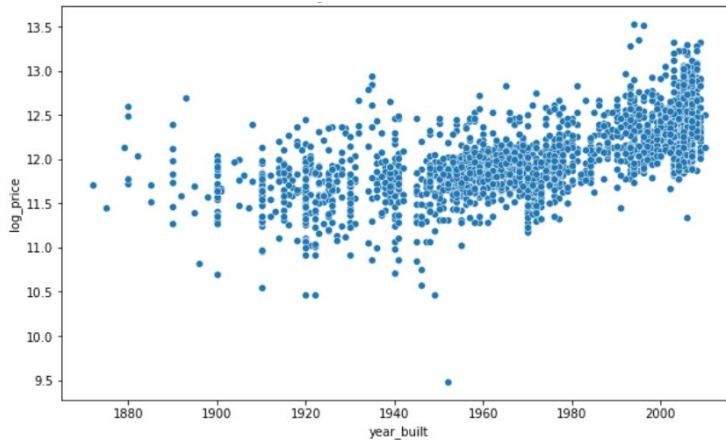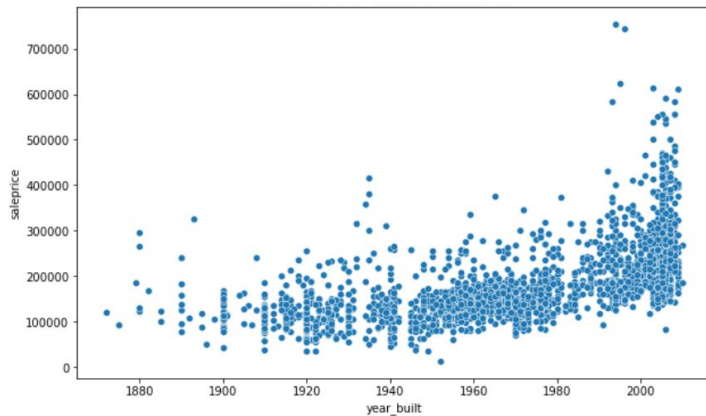
# Log Sale Price (Target)

**RMSE: 21,023**



Distribution of Sale Price

Skewness: 1.778

Grey Line: Mean

Green Line: Median

Distribution of Log Sale Price

Skewness: -0.003

Sale Price VS Log Sale Price

# Interaction Terms
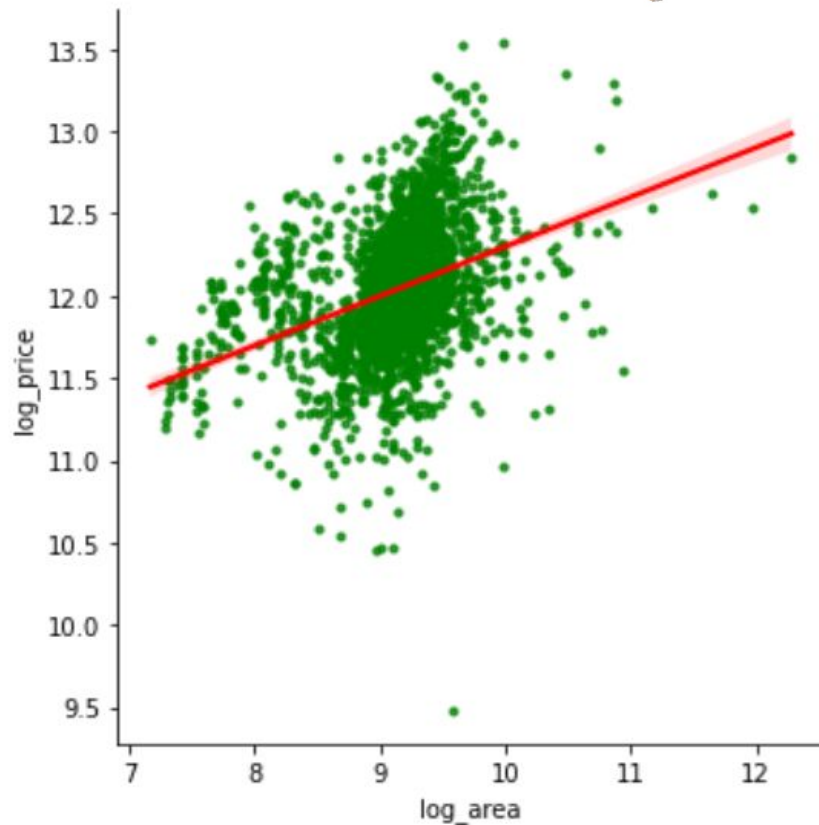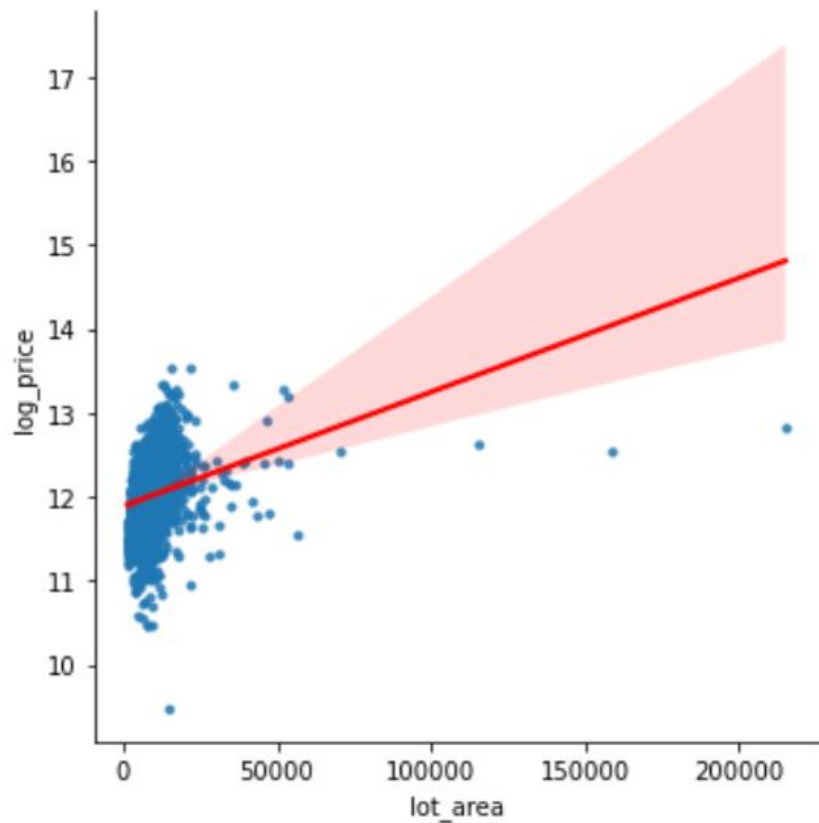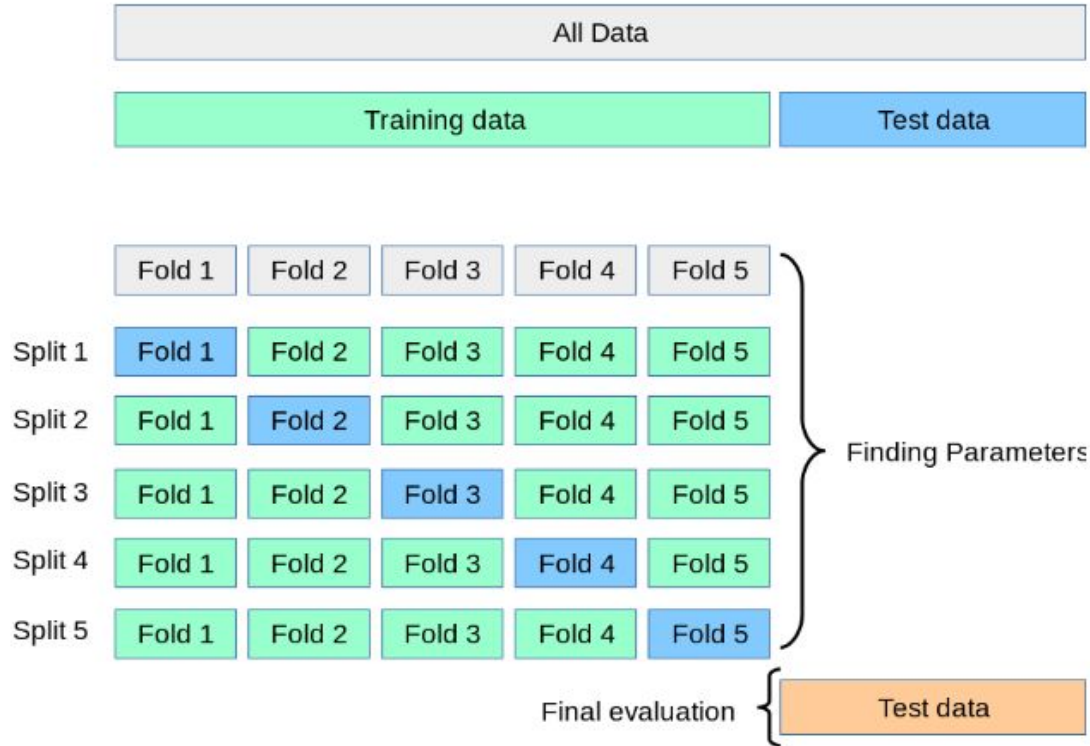


**Every Good House Need a Good Fireplace!**

# Interaction Terms

# Log feature

# More Power!



Use all data to training

RMSE: 23,351

The data says we need more data.

someecards
user card

# Summary

**Top 15 features**

Correlation with sale price

**01**

**Data Cleaning**

Remove outliers

**02**

**Pattern identification**

Group Dummy Variables

**03**

**Relationship between X&Y**

Log Transformation

**04**

**Thought Process**