

OVERVIEW

BUSINESS AND DATA UNDERSTANDING

CONTENTS

MODELING

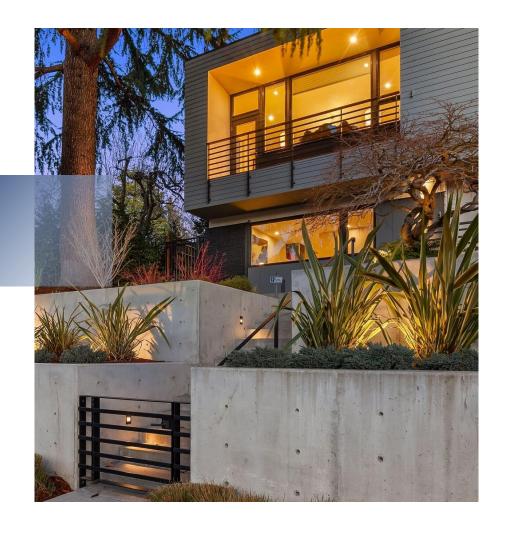
Present baseline and multiple linear regression models

REGRESSION RESULTS

Analyze R² of multiple regression model and assumptions of regression

RECOMMENDATION AND CONCLUSION





BUSINESS PROBLEM

- Real estate company purchases homes and tries to flip at premium
- Wants to use data on house sales to see if the purchase and sale price are profitable



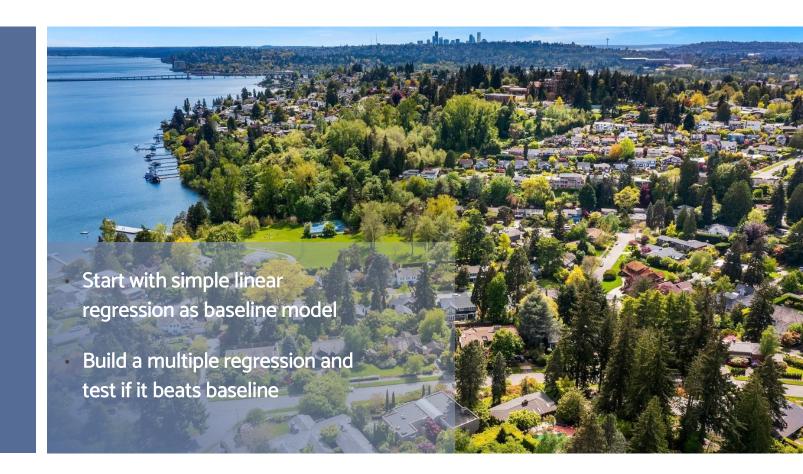
30,311

25

Non-null entries

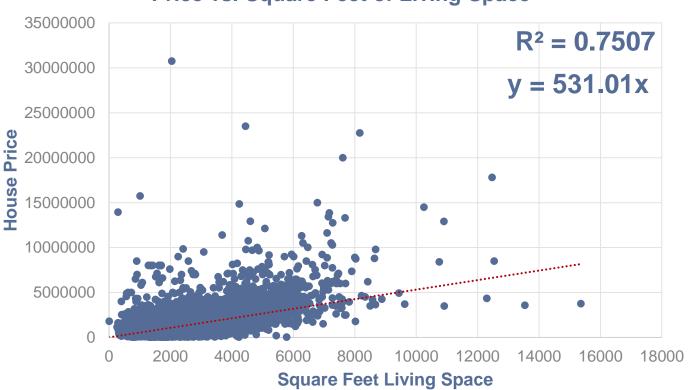
Columns; 23 aside from price and id to serve as predictors 0.61

Positive correlation between sqft_living and price



Baseline Model

Price vs. Square Feet of Living Space



 $y = -60,300x_1 + 102x_2 + 0.06x_2^2 - 89,580x_3 - 242x_4 + 169x_5 + 177x_6 + 207,100x_7 + 313,400x_8$

Variable	Column Name	Coefficient Value
x_1	Bedrooms	-60,300
x_2	Sqft Living Space	102
x_2^2	Sqft Living Space Squared	0.06
x_3	Floors	-89,580
x_4	Sqft Garage	-242
x_5	Sqft Patio	169
x_6	Year Built	177
x_7	Grade	207,100
x_8	View	313,400

$$y = (-60,300x_1) + 102x_2 + 0.06x_2^2 - 89,580x_3 - 242x_4 + 169x_5 + 177x_6 + 207,100x_7 + 313,400x_8$$

Variable	Column Name	Coefficient Value
x_1	Bedrooms	-60,300
x_2	Sqft Living Space	102
x_2^2	Sqft Living Space Squared	0.06
x_3	Floors	-89,580
<i>x</i> ₄	Sqft Garage	-242
x_5	Sqft Patio	169
x_6	Year Built	177
<i>x</i> ₇	Grade	207,100
x_8	View	313,400

$$y = -60,300x_1 + 102x_2 + 0.06x_2^2 - 89,580x_3 - 242x_4 + 169x_5 + 177x_6 + 207,100x_7 + 313,400x_8$$

Variable	Column Name	Coefficient Value
x_1	Bedrooms	-60,300
x_2	Sqft Living Space	102
x_{2}^{2}	Sqft Living Space Squared	0.06
<i>x</i> ₃	Floors	-89,580
<i>x</i> ₄	Sqft Garage	-242
x_5	Sqft Patio	169
x_6	Year Built	177
x_7	Grade	207,100
<i>x</i> ₈	View	313,400

$$y = -60,300x_1 + 102x_2 + 0.06x_2^2 - 89,580x_3 - 242x_4 + 169x_5 + 177x_6 + 207,100x_7 + 313,400x_8$$

Variable	Column Name	Coefficient Value
x_1	Bedrooms	-60,300
x_2	Sqft Living Space	102
x_2^2	Sqft Living Space Squared	0.06
<i>x</i> ₃	Floors	-89,580
<i>x</i> ₄	Sqft Garage	-242
x_5	Sqft Patio	169
<i>x</i> ₆	Year Built	177
x_7	Grade	207,100
x_8	View	313,400

$$y = -60,300x_1 + 102x_2 + 0.06x_2^2 - 89,580x_3 - 242x_4 + 169x_5 + 177x_6 + 207,100x_7 + 313,400x_8$$

Variable	Column Name	Coefficient Value
x_1	Bedrooms	-60,300
x_2	Sqft Living Space	102
x_2^2	Sqft Living Space Squared	0.06
x_3	Floors	-89,580
<i>x</i> ₄	Sqft Garage	-242
x_5	Sqft Patio	169
x_6	Year Built	177
<i>x</i> ₇	Grade	207,100
<i>x</i> ₈	View	313,400

$$y = -60,300x_1 + 102x_2 + 0.06x_2^2 - 89,580x_3 - 242x_4 + 169x_5 + 177x_6 + 207,100x_7 + 313,400x_8$$

Variable	Column Name	Coefficient Value
x_1	Bedrooms	-60,300
x_2	Sqft Living Space	102
x_2^2	Sqft Living Space Squared	0.06
x_3	Floors	-89,580
<i>x</i> ₄	Sqft Garage	-242
<i>x</i> ₅	Sqft Patio	169
x_6	Year Built	177
x_7	Grade	207,100
<i>x</i> ₈	View	313,400

$$y = -60,300x_1 + 102x_2 + 0.06x_2^2 - 89,580x_3 - 242x_4 + 169x_5 + 177x_6 + 207,100x_7 + 313,400x_8$$

Variable	Column Name	Coefficient Value
x_1	Bedrooms	-60,300
x_2	Sqft Living Space	102
x_2^2	Sqft Living Space Squared	0.06
x_3	Floors	-89,580
<i>x</i> ₄	Sqft Garage	-242
x_5	Sqft Patio	169
x_6	Year Built	177
<i>x</i> ₇	Grade	207,100
<i>x</i> ₈	View	313,400

$$y = -60,300x_1 + 102x_2 + 0.06x_2^2 - 89,580x_3 - 242x_4 + 169x_5 + 177x_6 + 207,100x_7 + 313,400x_8$$

Variable	Column Name	Coefficient Value
x_1	Bedrooms	-60,300
x_2	Sqft Living Space	102
x_2^2	Sqft Living Space Squared	0.06
x_3	Floors	-89,580
<i>x</i> ₄	Sqft Garage	-242
x_5	Sqft Patio	169
x_6	Year Built	177
<i>x</i> ₇	Grade	207,100
x_8	View	313,400

$$y = -60,300x_1 + 102x_2 + 0.06x_2^2 - 89,580x_3 - 242x_4 + 169x_5 + 177x_6 + 207,100x_7 + 313,400x_8$$

Variable	Column Name	Coefficient Value
x_1	Bedrooms	-60,300
x_2	Sqft Living Space	102
x_2^2	Sqft Living Space Squared	0.06
x_3	Floors	-89,580
x ₄	Sqft Garage	-242
x_5	Sqft Patio	169
x_6	Year Built	177
x_7	Grade	207,100
x_8	View	313,400



 $y = -60,300x_1 + 102x_2 + 0.06x_2^2 - 89,580x_3 - 242x_4 + 169x_5 + 177x_6 + 207,100x_7 + 313,400x_8$

$R^2 = 0.78$

Our model explains 78% of total variance in house price and all coefficients are statistically significant

LINEAR

Unfortunately, statistical testing showed that our model was not linear

NORMALLY DISTRIBUTED RESIDUALS

Statistical testing showed us that our residuals were in fact normally distributed

INDEPENDENCE OF VARIABLES

In our modeling process, we worked hard to elimante all variables that were highly correlated

HOMOSKEDASTICITY

The scatter of the errors (difference between actual and predicted values) is even

05

USE OUR MODEL, BUT CAUTIOUSLY

- Our model picks up on a decent chunk of variance in home prices
- It is statistically significant
- However it isn't linear
- Good practice to use intuitive real estate methods and use model to check



