

# Mod 1 Project – King County Housing Data

Thomas O’Gara

Part time data science

Instructor: Victor Geislinger

# Problem Statement:

- How do we go about predicting home values in King County?



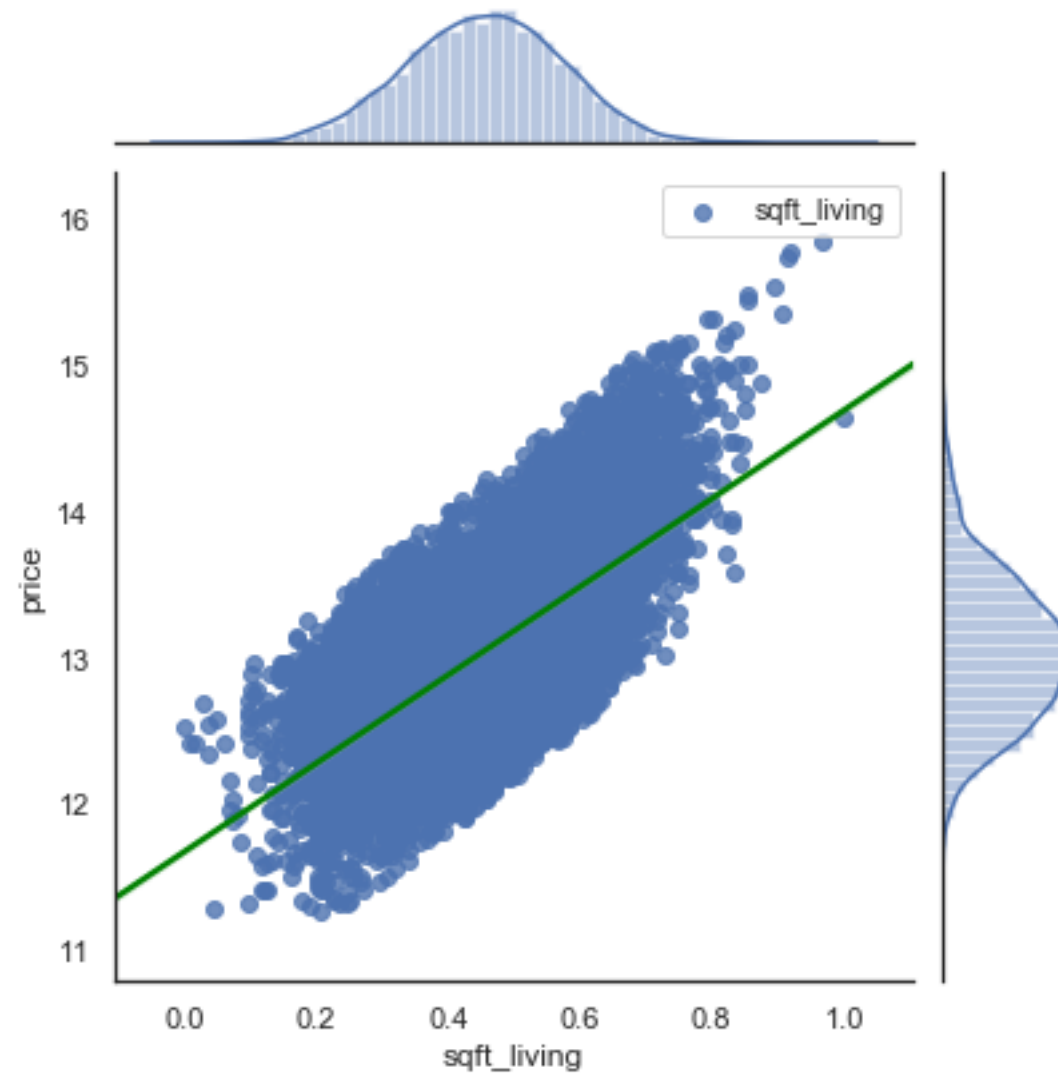
# Methodology – “OSEMN”

- **O**btained housing data on 21,000+ King County homes
- **S**crubbed the data for missing values, non-numbers, placeholder values
- **E**xploratory data analysis (EDA)
- **M**odeled the data using multiple linear regression
- **I**nterpreted the data to determine the most important factors in valuing a home in King County

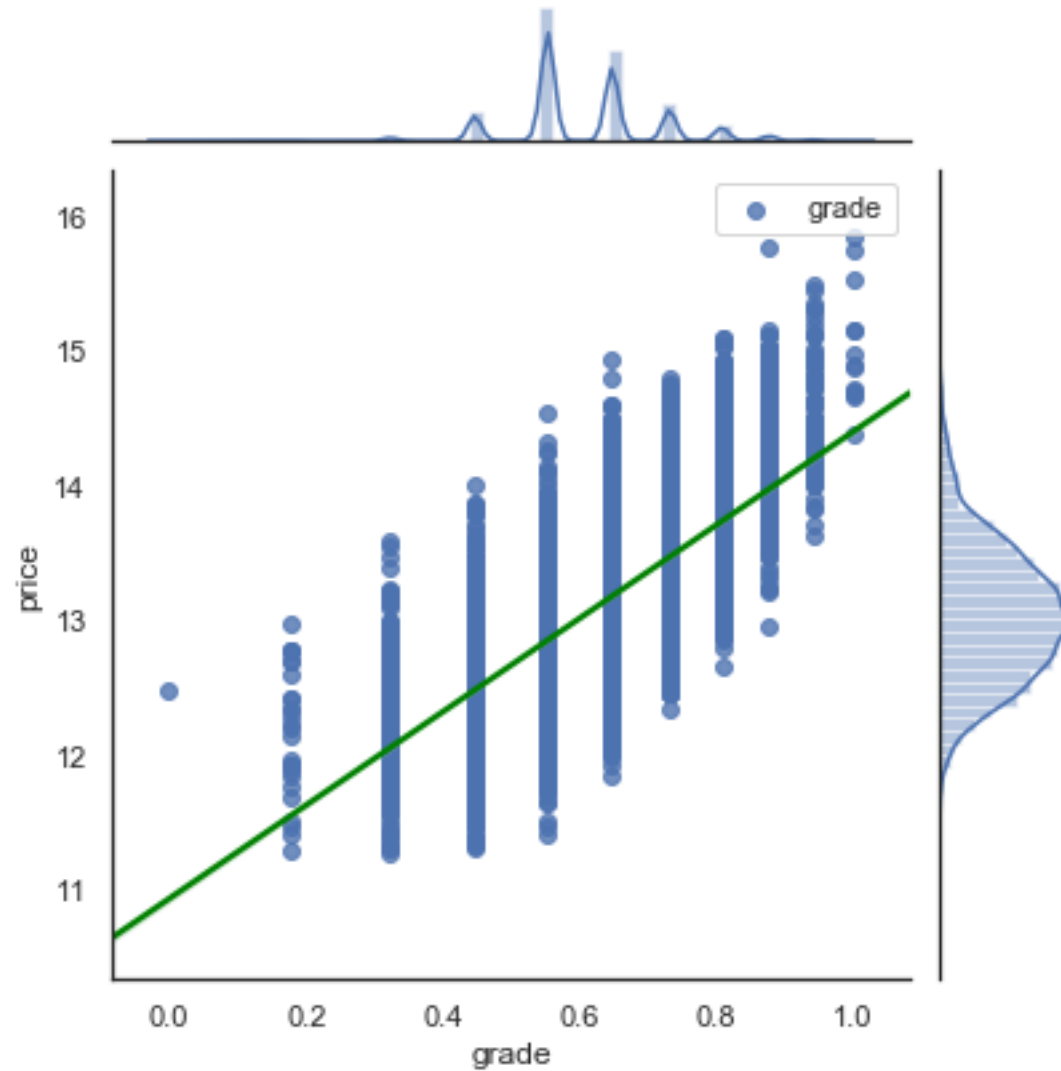
# Most important variables to predict home prices

1. Square footage of the home
2. Grade/build quality
3. Square footage of the lot
4. Square footage of neighbor's homes
5. Location/neighborhood

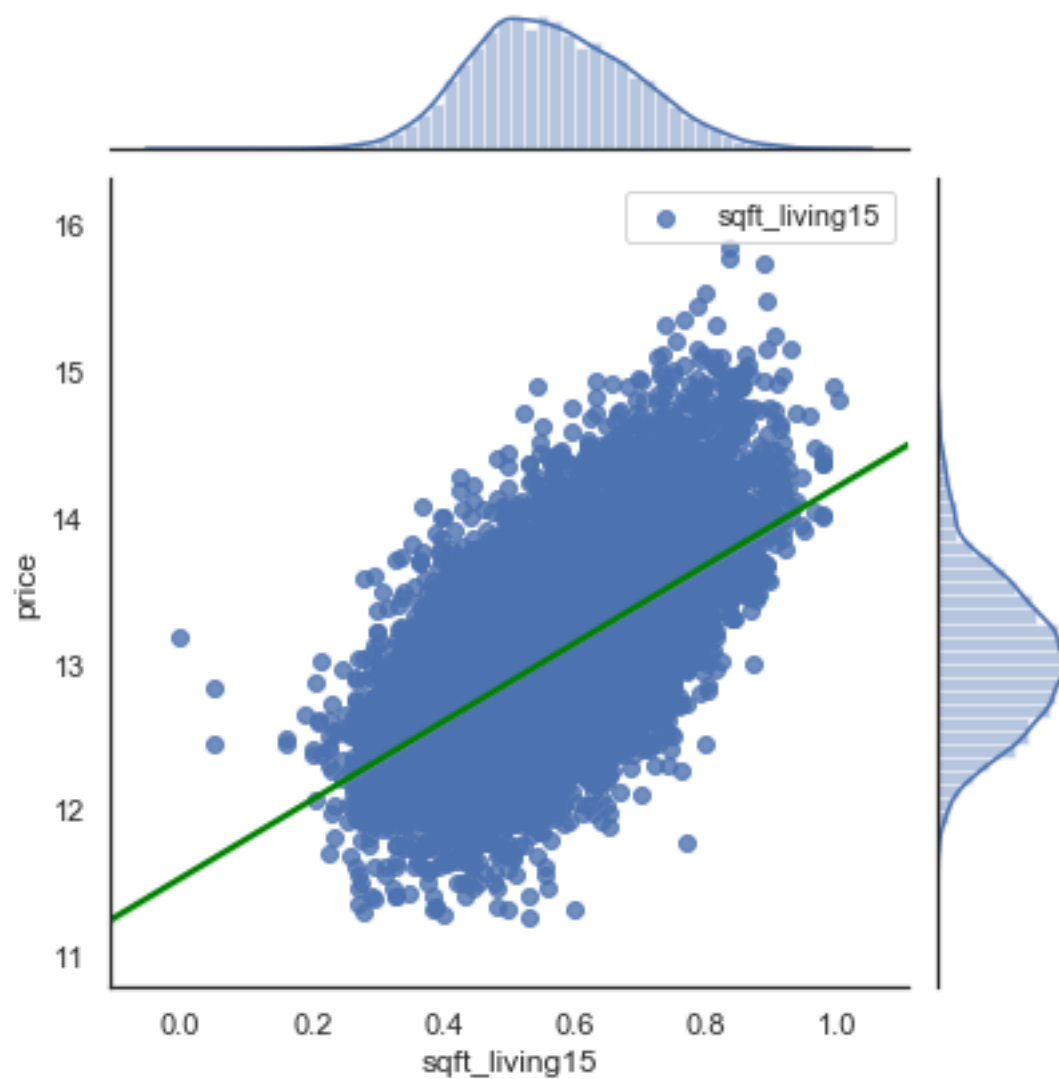
# Square footage of the home



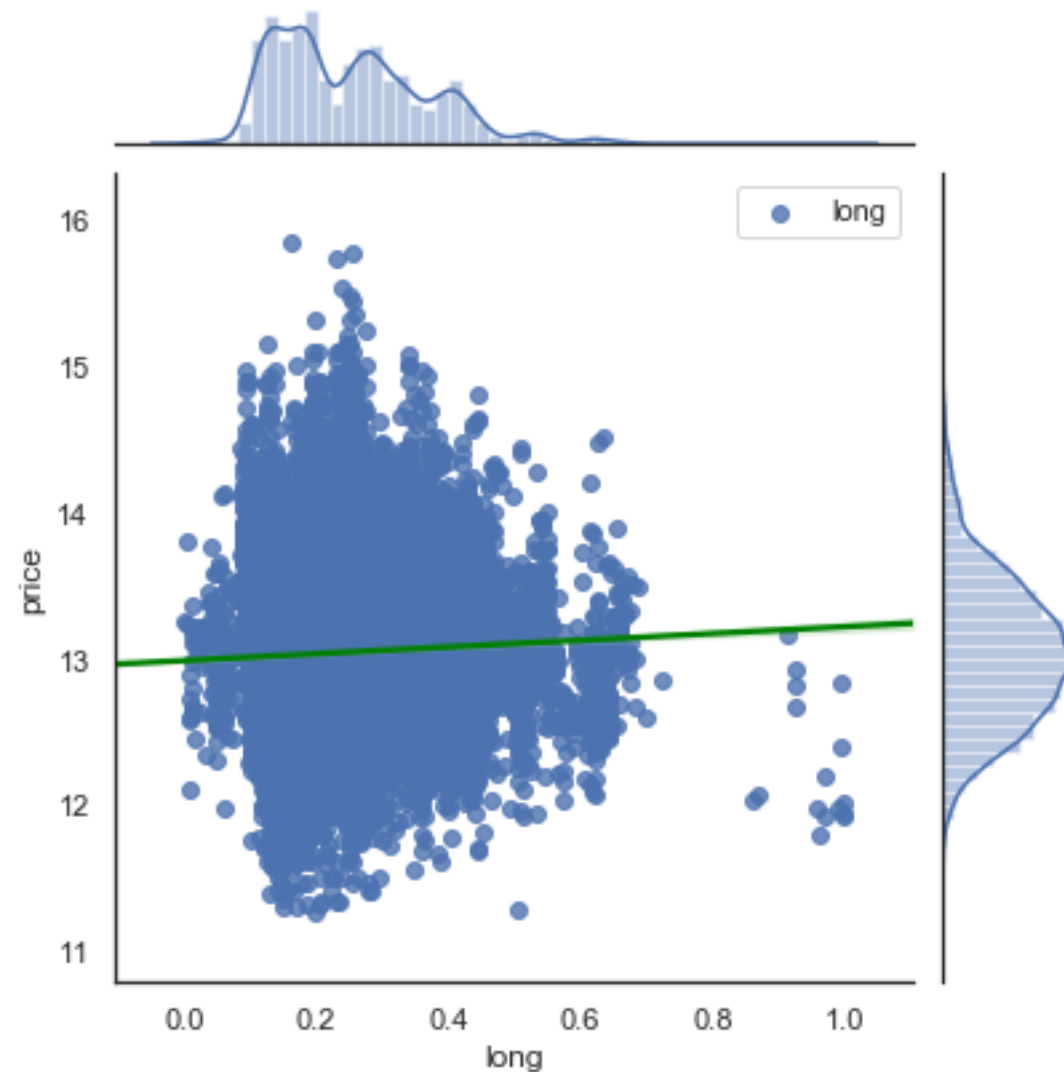
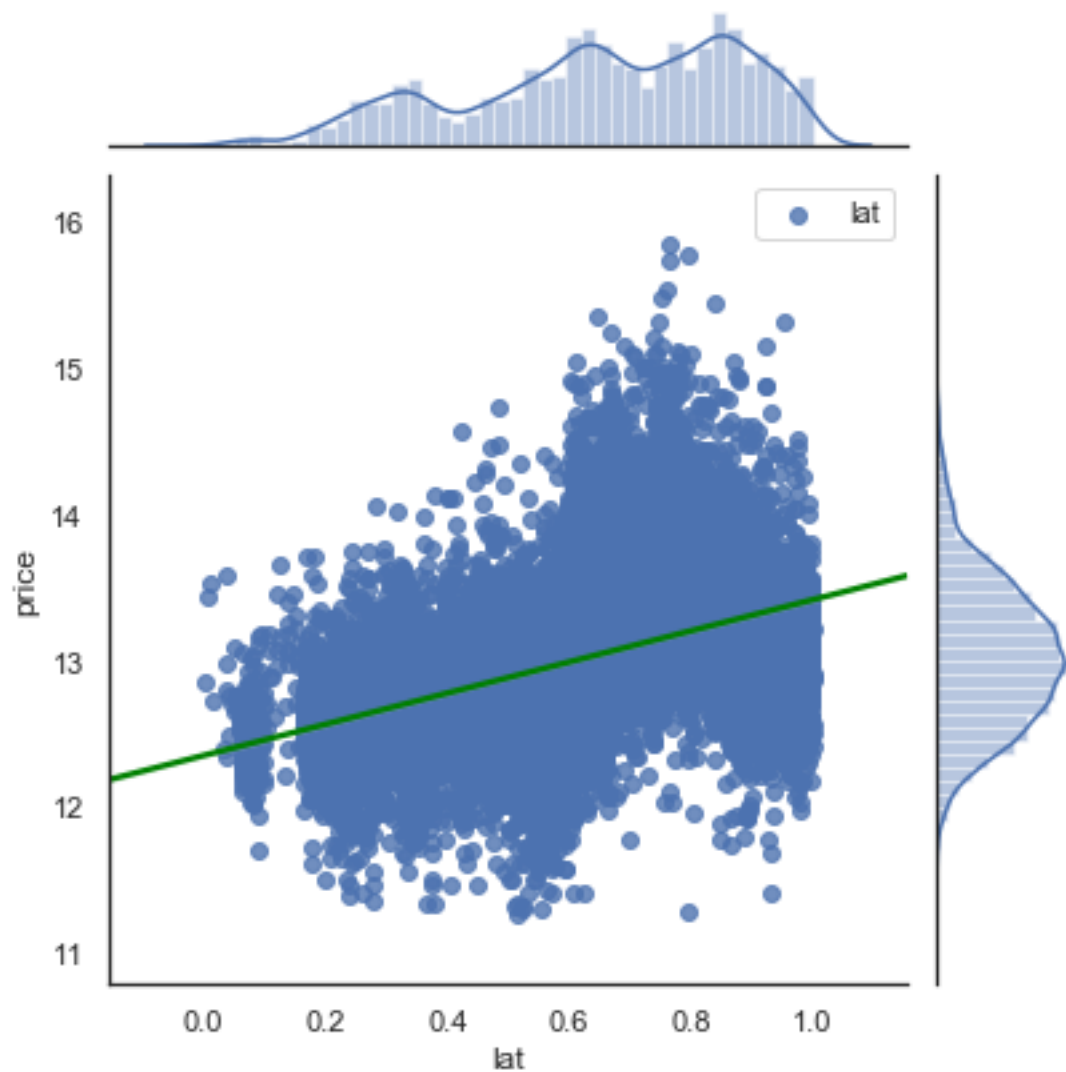
# Grade/build quality



# Square footage of neighbor's homes

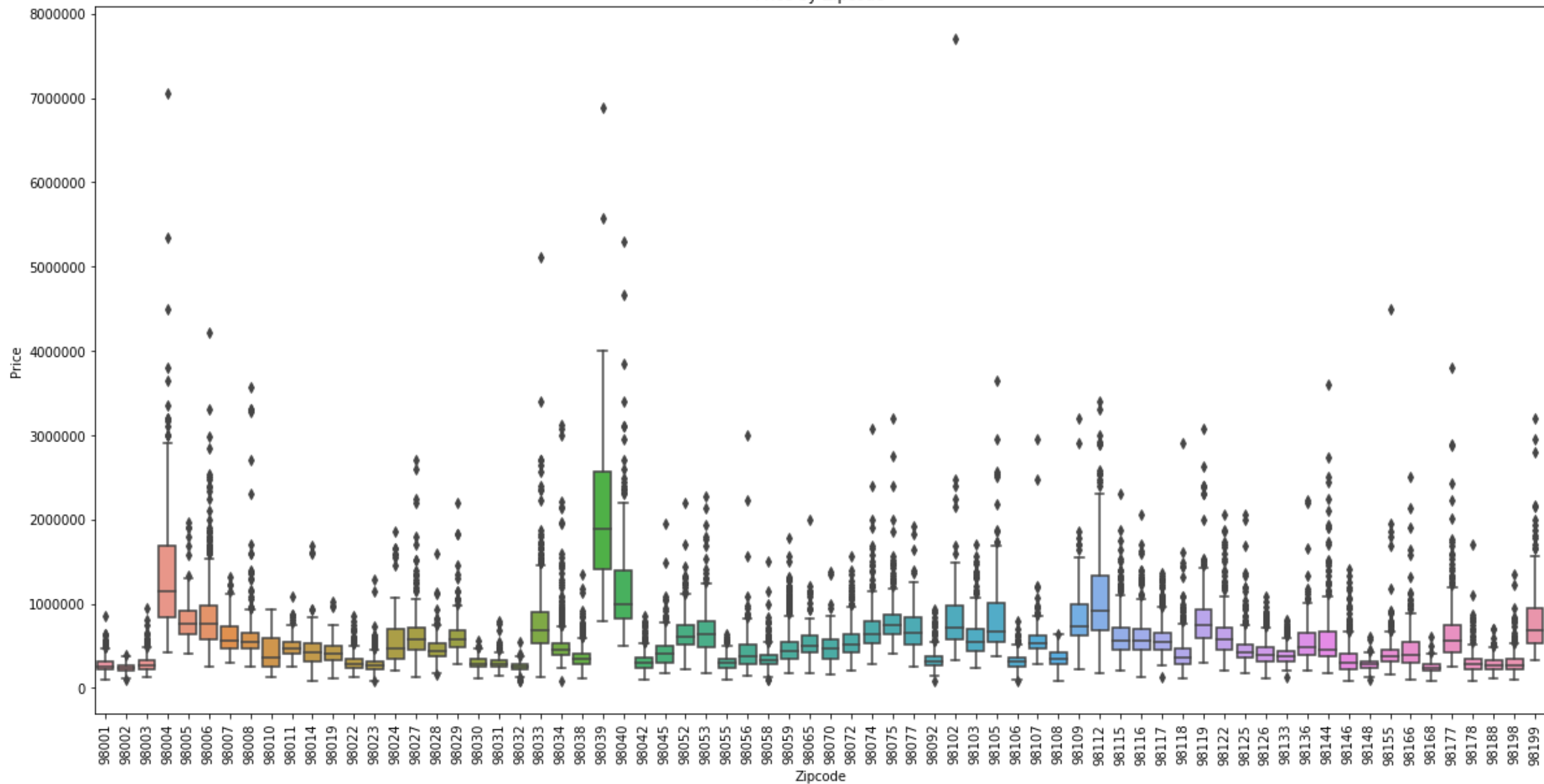


# Location - where do you want to be?





Price by Zipcode

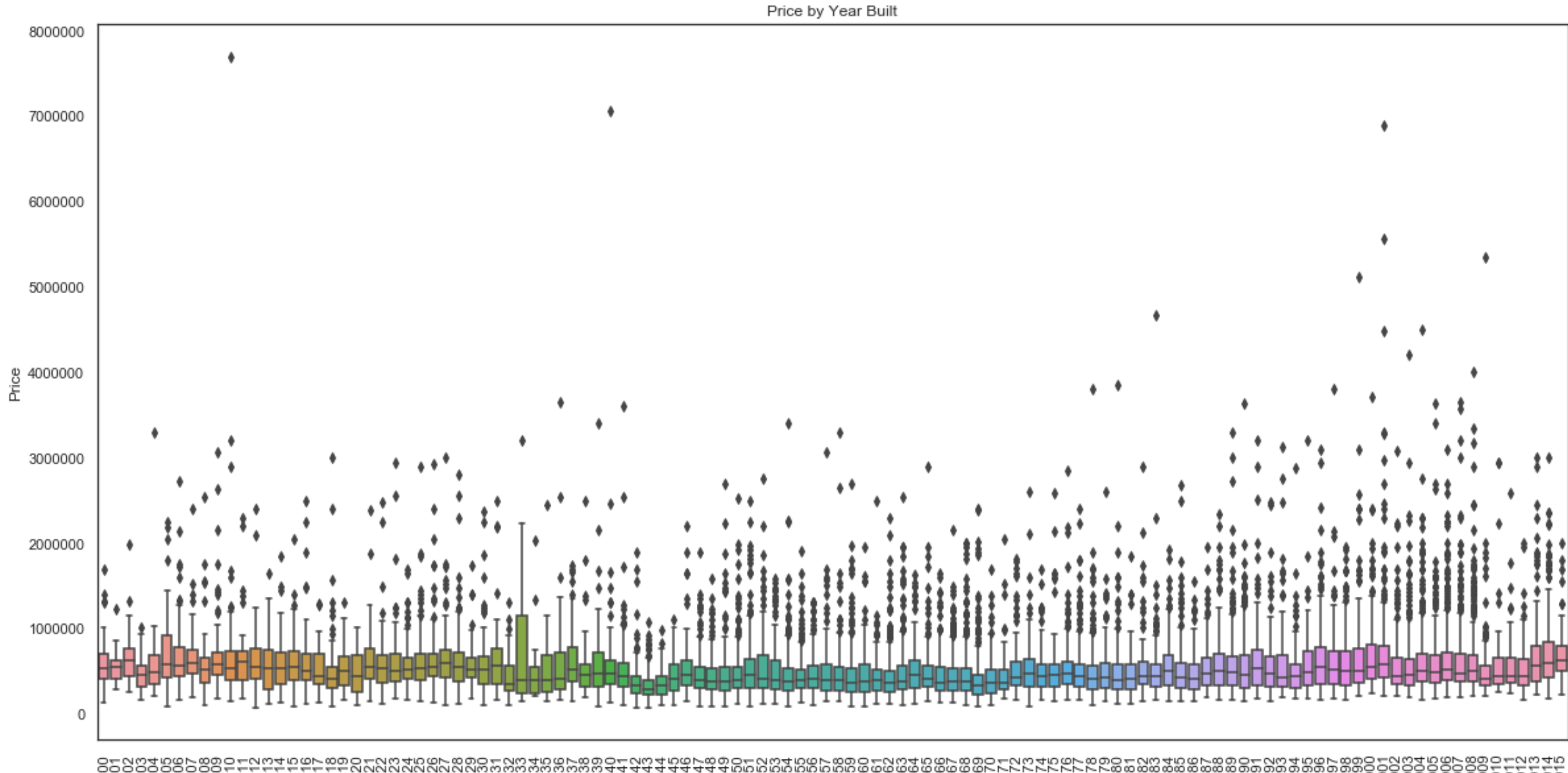


# Location continued:

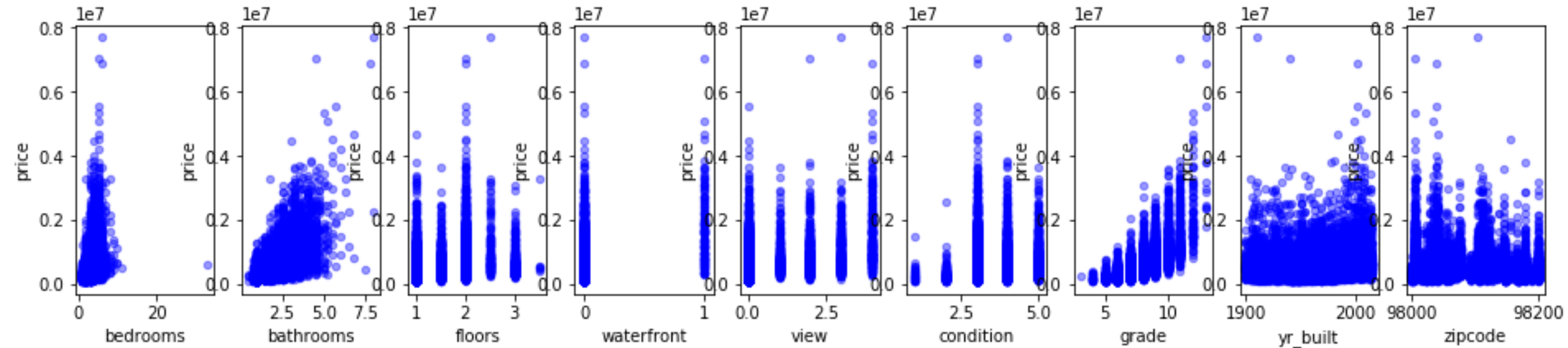
## Most expensive zip codes:

- 98004: Medina
- 98039: Bellevue
- 98040: Mercer Island
- 98112: Northeast Seattle
- 98199: Northwest Seattle

# Is a newer home more valuable?



# What about other variables?



# Model Overview

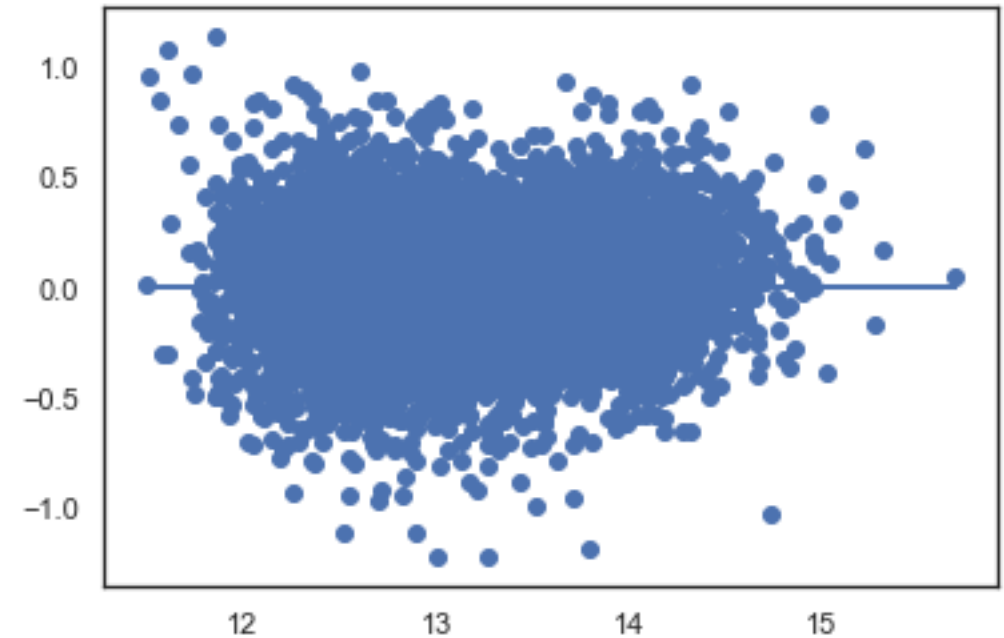
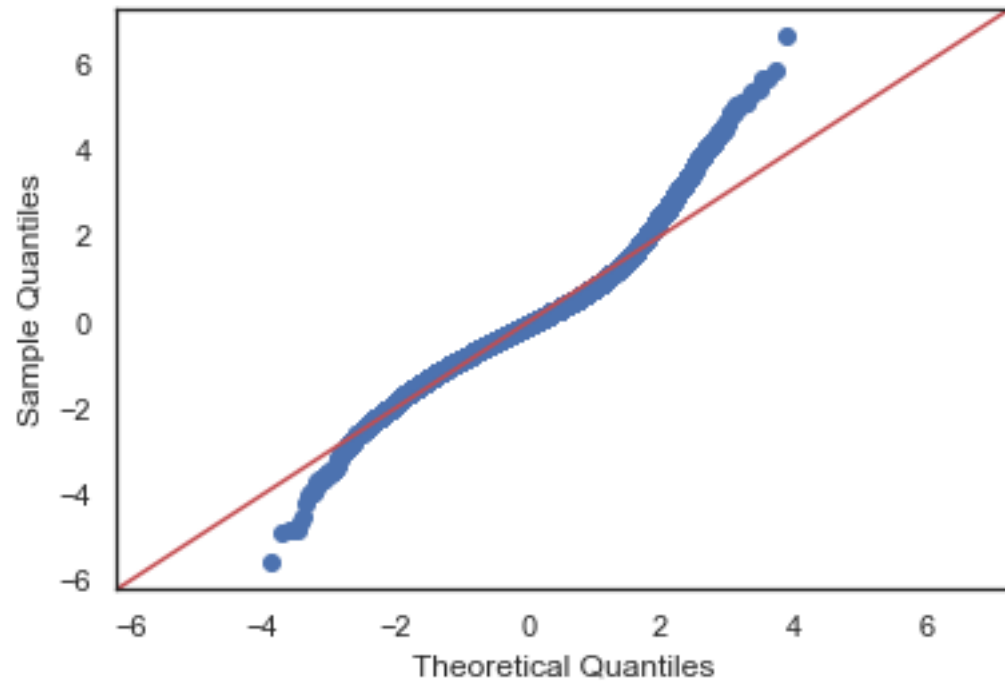
## OLS Regression Results

<b>Dep. Variable:</b>	price	<b>R-squared:</b>	0.878
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.878
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	1646.
<b>Date:</b>	Wed, 12 Jun 2019	<b>Prob (F-statistic):</b>	0.00
<b>Time:</b>	10:38:44	<b>Log-Likelihood:</b>	5811.6
<b>No. Observations:</b>	21082	<b>AIC:</b>	-1.144e+04
<b>Df Residuals:</b>	20989	<b>BIC:</b>	-1.070e+04
<b>Df Model:</b>	92		
<b>Covariance Type:</b>	nonrobust		

<b>Omnibus:</b>	1302.533	<b>Durbin-Watson:</b>	2.001
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	5939.522
<b>Skew:</b>	-0.055	<b>Prob(JB):</b>	0.00
<b>Kurtosis:</b>	5.598	<b>Cond. No.</b>	281.

	coef	std err	t	P> t	[0.025	0.975]
<b>Intercept</b>	10.9549	0.059	186.655	0.000	10.840	11.070
<b>bedrooms</b>	-0.2552	0.024	-10.684	0.000	-0.302	-0.208
<b>bathrooms</b>	0.1798	0.016	10.950	0.000	0.148	0.212
<b>sqft_living</b>	1.5118	0.027	55.100	0.000	1.458	1.566
<b>sqft_lot</b>	0.6256	0.030	20.897	0.000	0.567	0.684
<b>grade</b>	1.0027	0.022	46.034	0.000	0.960	1.045
<b>sqft_basement</b>	-0.0624	0.005	-13.359	0.000	-0.072	-0.053
<b>lat</b>	0.3064	0.046	6.703	0.000	0.217	0.396
<b>long</b>	-0.5190	0.064	-8.059	0.000	-0.645	-0.393
<b>sqft_living15</b>	0.4074	0.019	21.551	0.000	0.370	0.444
<b>sqft_lot15</b>	-0.0980	0.030	-3.290	0.001	-0.156	-0.040
<b>floor[0]</b>	0.0244	0.005	4.806	0.000	0.014	0.034
<b>floor[1]</b>	0.0085	0.004	2.014	0.044	0.000	0.017
<b>floor[2]</b>	0.0504	0.016	3.235	0.001	0.020	0.081
<b>floor[3]</b>	-0.0703	0.009	-7.446	0.000	-0.089	-0.052
<b>floor[4]</b>	-0.0362	0.070	-0.518	0.605	-0.173	0.101
<b>view[0]</b>	0.1188	0.011	11.285	0.000	0.098	0.139
<b>view[1]</b>	0.1151	0.006	17.798	0.000	0.102	0.128
<b>view[2]</b>	0.2029	0.009	23.130	0.000	0.186	0.220
<b>view[3]</b>	0.4891	0.011	44.462	0.000	0.468	0.511
<b>condition_2</b>	0.1269	0.038	3.360	0.001	0.053	0.201
<b>condition_3</b>	0.2417	0.035	6.880	0.000	0.173	0.311
<b>condition_4</b>	0.2755	0.035	7.838	0.000	0.207	0.344
<b>condition_5</b>	0.3416	0.035	9.658	0.000	0.272	0.411

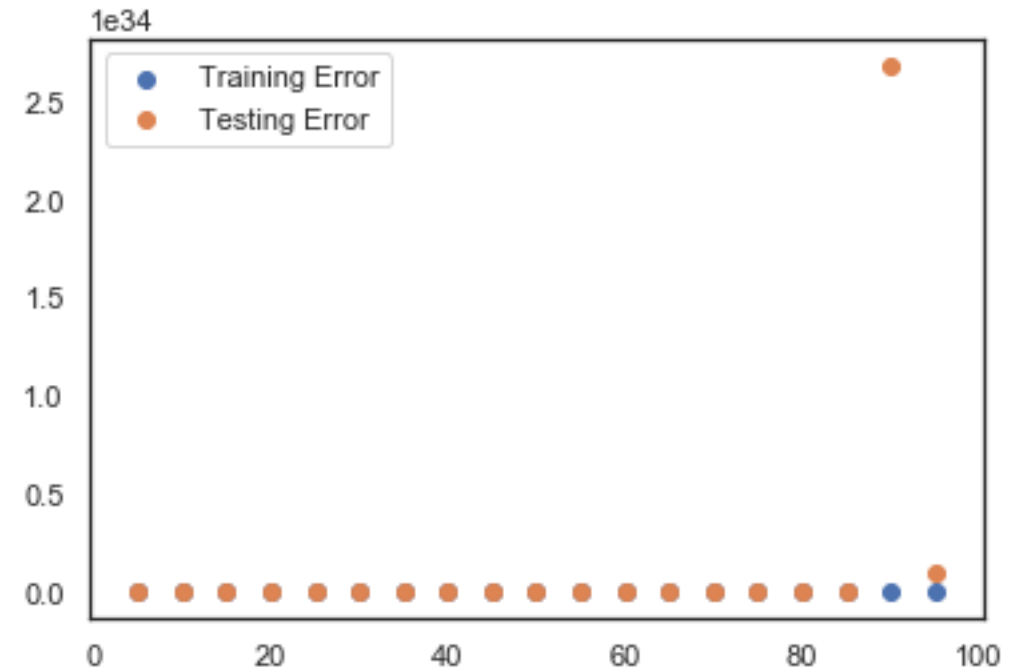
# Normality and Homoscedasticity assumptions



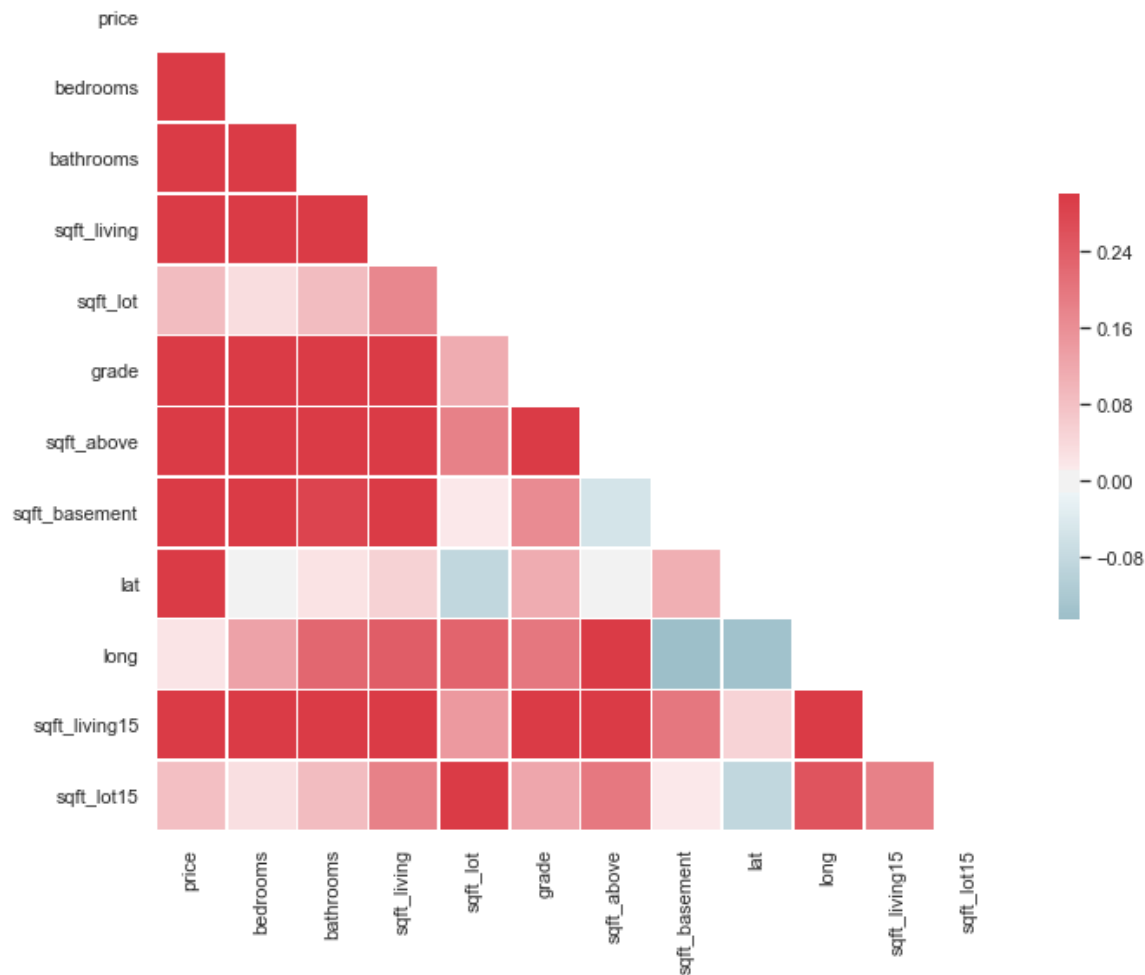
# Is this model a good predictor of future home prices?

```
mse_train = np.sum((y_train-y_hat_train)**2)/len(y_train)
mse_test = np.sum((y_test-y_hat_test)**2)/len(y_test)
print('Train Mean Squared Error:', mse_train)
print('Test Mean Squared Error:', mse_test)
```

Train Mean Squared Error: 0.034076607061633214  
Test Mean Squared Error: 0.032997092011984586



# What are potential problems with this model?



```
#Checking for Multicollinearity using variance inflation factor.  
#this accounts for multicollinearity with a relation of 3 or more variables v  
from statsmodels.stats.outliers_influence import variance_inflation_factor
```

```
X = subset.drop(['price'], axis=1)  
vif = [variance_inflation_factor(X.values, i) for i in range(X.shape[1])]  
list(zip(X.columns, vif))
```

```
[('bedrooms', 27.43936322112261),  
 ('bathrooms', 28.899237707052393),  
 ('sqft_living', 102.81077627409047),  
 ('sqft_lot', 22.845588831113602),  
 ('grade', 112.24690739331463),  
 ('sqft_basement', 3.1025613421357163),  
 ('lat', 273.2524962797792),  
 ('long', 193.36742543893328),  
 ('sqft_living15', 72.38936688954992),  
 ('floor', 1.4059878707820597),  
 ('floor', 4.266375251366034),  
 ('floor', 1.1034480605684647),  
 ('floor', 1.5573844297385897),  
 ('floor', 1.0107382278913626),  
 ('view', 1.0703526808346429),  
 ('view', 1.1481992787171456),  
 ('view', 1.1259385931881711),  
 ('view', 1.1098055988207716),  
 ('condition_2', 4.837279868308839),  
 ('condition_3', 321.7072792016948),  
 ('condition_4', 129.5277054533042),  
 ('condition_5', 39.66082642763849),
```



# Business recommendation for home values

- Square footage of living space is key
- High grade building
  - Doesn't need to be a new build – just needs to be high quality workmanship
- Pick a neighborhood with large homes nearby
  - Better to be the worst house in a good neighborhood
- Northwest area of King County is best
  - Focus on top 5 zip codes
  - Close to the water

Thank you!

Questions?