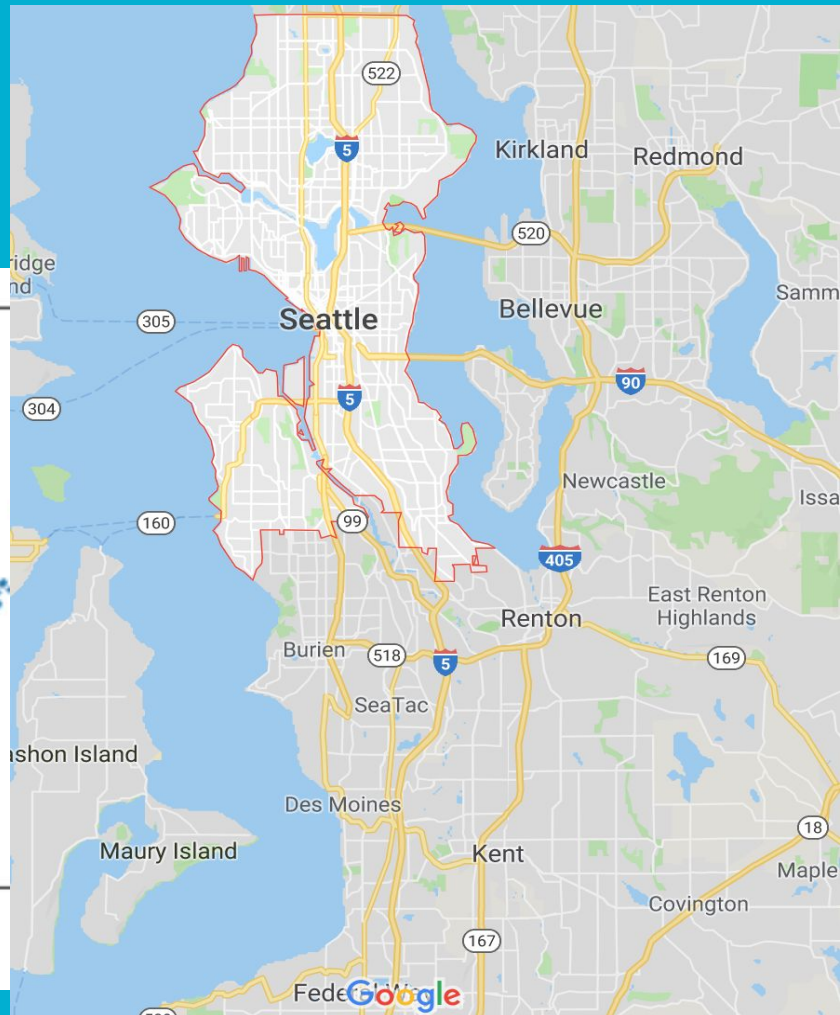
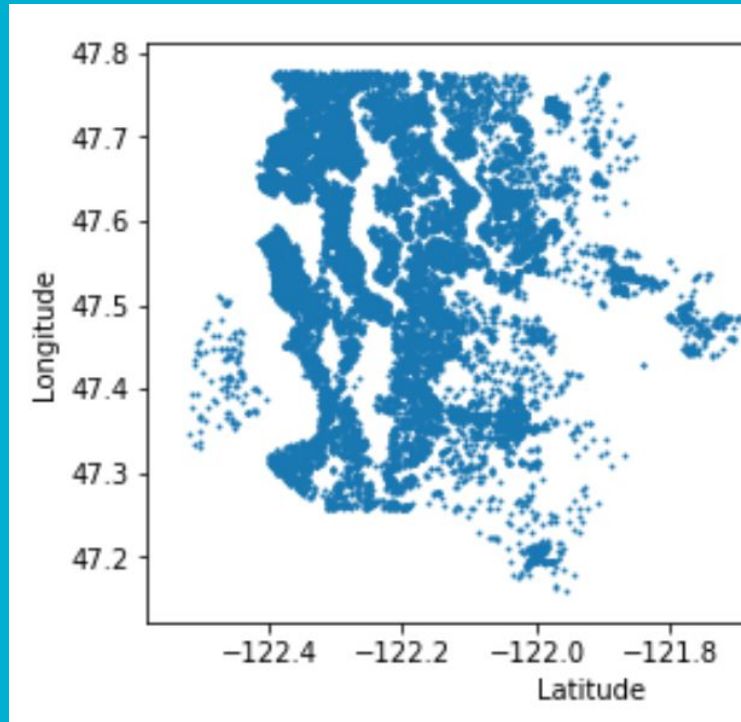


# Predicting Home Sale Prices in Greater Seattle

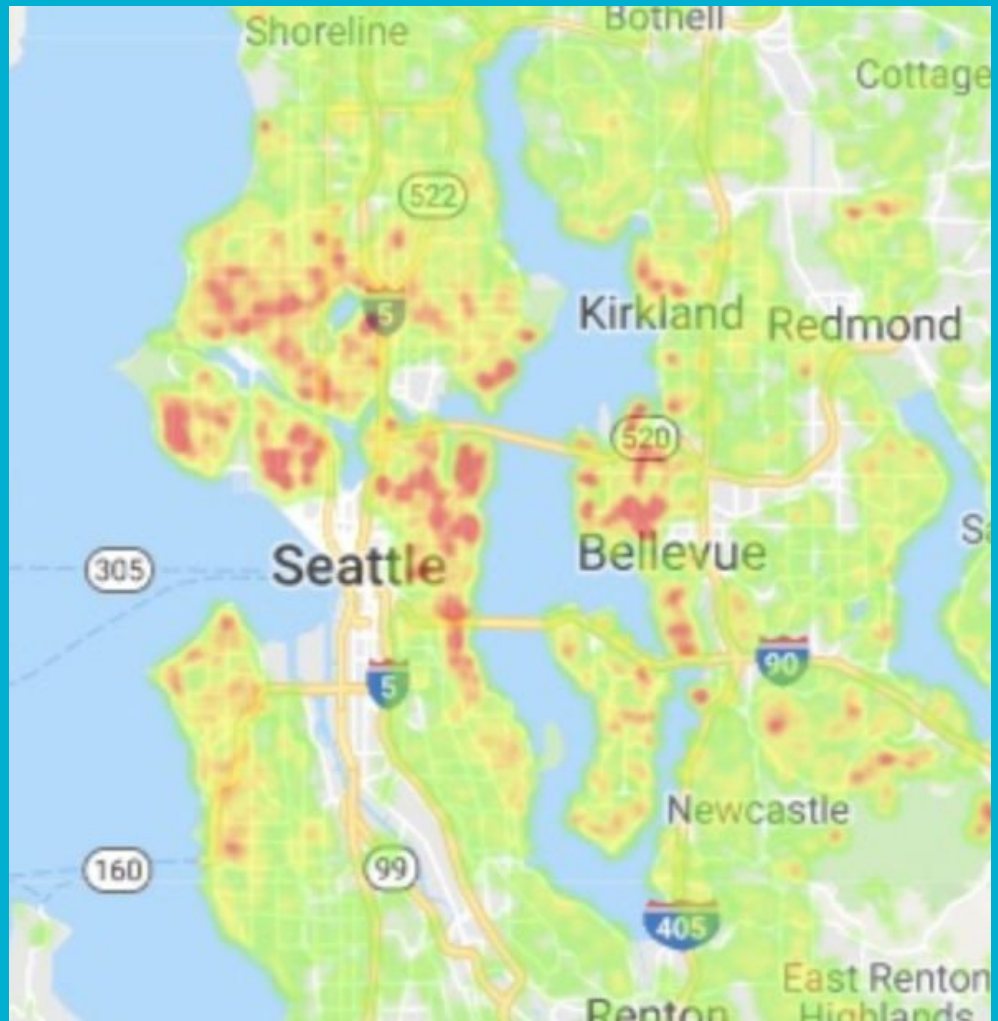
---

Leveraging a King County House Sales dataset for insights

# Scatter Plot of Latitude & Longitude



# Heat Map Highlighting Housing Density



# King County House Sales Dataset

---

## Original Dataset (Sales from 2014-2015)

- 21,597 instances of home sales (observations, rows)
- 21 columns or features

## Cleaned Dataset

- 21,313 observations after removing 284 outliers (or extreme influencers)
- Maintained 98.7% of the data set (trimmed only 1.3% of the data)
- Included 10 of the original 21 features in our model as predictors of price

# Model Description

---

Developed a model with 85% accuracy

- Important predictors...
  - Square footage of the home
  - Location
  - Lot size
  - Cumulative impact of condition, year built & renovation
- Outcome (or target) variable → log of price

# Training Data

## 75%

```
In [45]: model = ols(formula= formula, data=train).fit()  
model.summary()
```

Out[45]: OLS Regression Results

<b>Dep. Variable:</b>	log_price	<b>R-squared:</b>	0.859			
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.857			
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	491.2			
<b>Date:</b>	Thu, 20 Jun 2019	<b>Prob (F-statistic):</b>	0.00			
<b>Time:</b>	10:37:31	<b>Log-Likelihood:</b>	3675.8			
<b>No. Observations:</b>	15984	<b>AIC:</b>	-6960.			
<b>Df Residuals:</b>	15788	<b>BIC:</b>	-5454.			
<b>Df Model:</b>	195					
<b>Covariance Type:</b>	nonrobust					
	<b>coef</b>	<b>std err</b>	<b>t</b>	<b>P&gt; t </b>	<b>[0.025</b>	<b>0.975]</b>
<b>Intercept</b>	41.8020	0.261	159.971	0.000	41.290	42.314

# Test Data

## 25%

```
In [46]: model = ols(formula= formula, data=test).fit()  
model.summary()
```

Out[46]:

OLS Regression Results

<b>Dep. Variable:</b>	log_price		<b>R-squared:</b>	0.865		
<b>Model:</b>	OLS		<b>Adj. R-squared:</b>	0.860		
<b>Method:</b>	Least Squares		<b>F-statistic:</b>	168.5		
<b>Date:</b>	Thu, 20 Jun 2019	<b>Prob (F-statistic):</b>	0.00			
<b>Time:</b>	10:42:10	<b>Log-Likelihood:</b>	1468.4			
<b>No. Observations:</b>	5329		<b>AIC:</b>	-2545.		
<b>Df Residuals:</b>	5133		<b>BIC:</b>	-1255.		
<b>Df Model:</b>	195					
<b>Covariance Type:</b>	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
<b>Intercept</b>	41.1340	0.440	93.442	0.000	40.271	41.997

# Conclusions

---

- Location, location, location  
[Next iteration of model will leverage latitude & longitude]
- Major on the majors  
[# of bathrooms, neighborhood & grade all captured by the proxy, square footage]



Questions?

# Multicollinearity

---

```
In [17]: df['sqft_living'].corr(df['sqft_above'])
```

```
Out[17]: 0.8764477590354981
```

```
In [18]: df['sqft_living'].corr(df['bathrooms'])
```

```
Out[18]: 0.7557576009502521
```

```
In [19]: df['sqft_living'].corr(df['grade'])
```

```
Out[19]: 0.7627790466721344
```

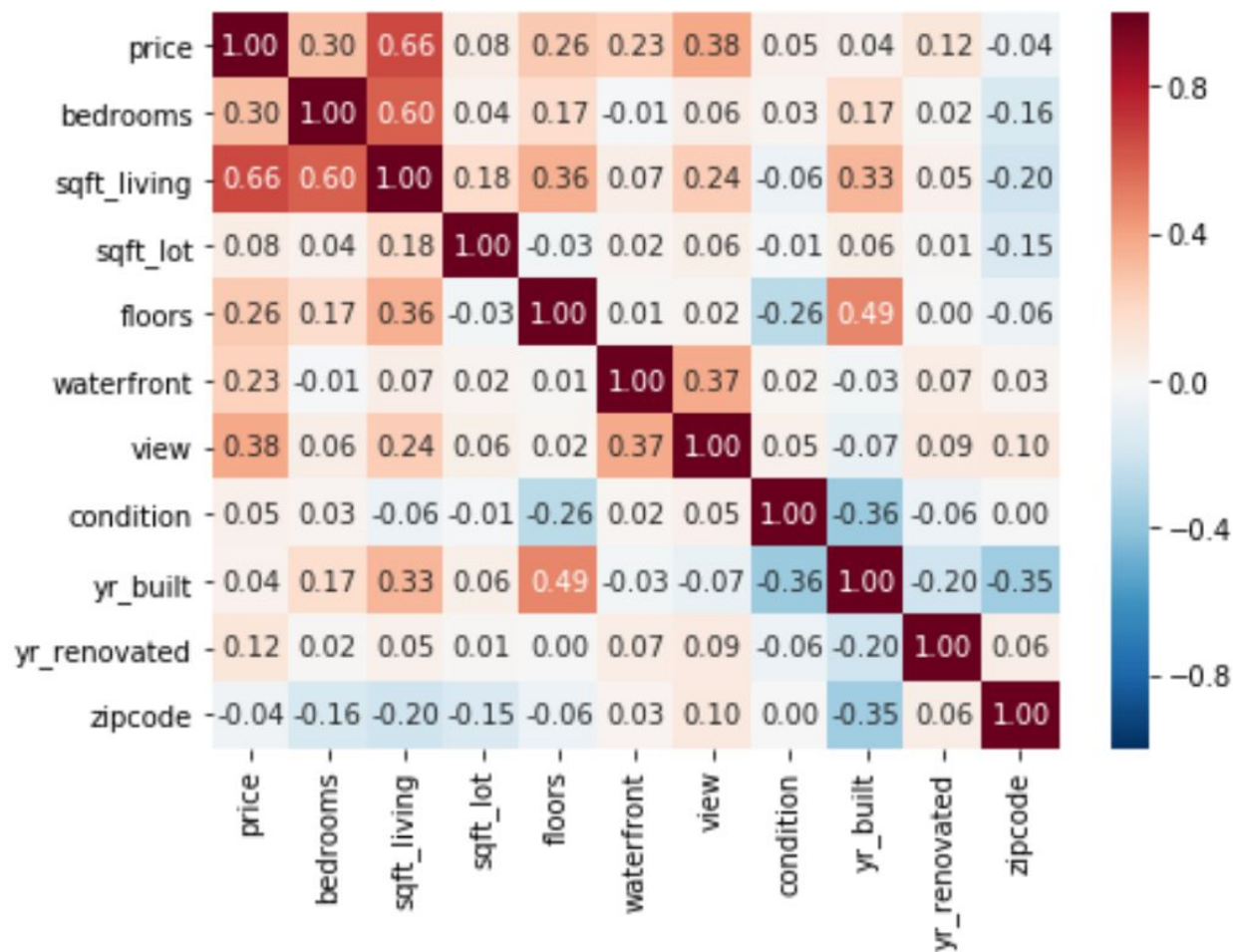
```
In [20]: df['sqft_living'].corr(df['sqft_living15'])
```

```
Out[20]: 0.7564015282475002
```

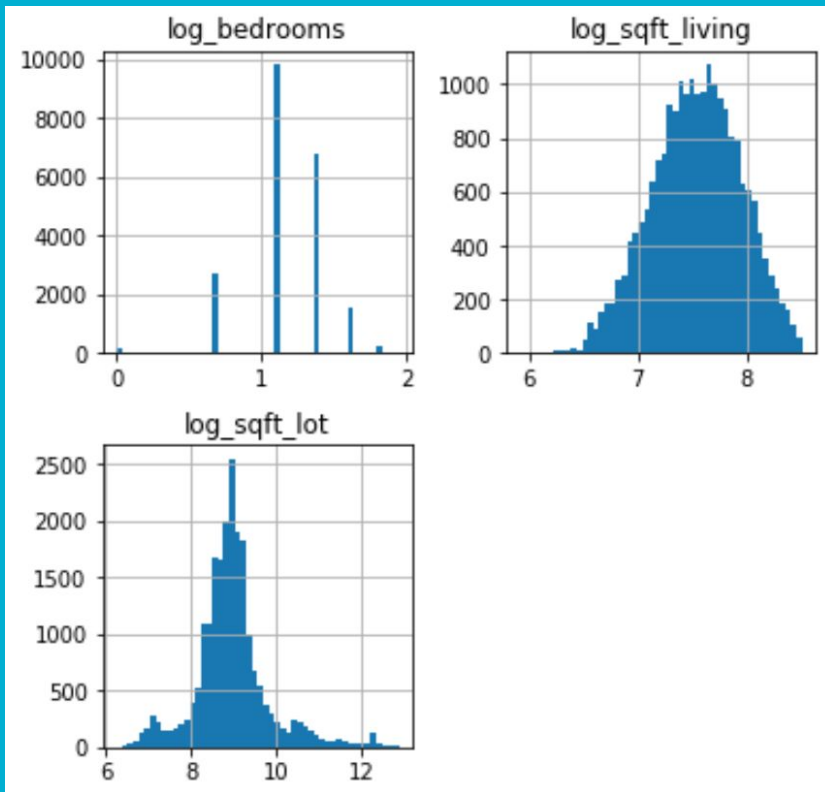
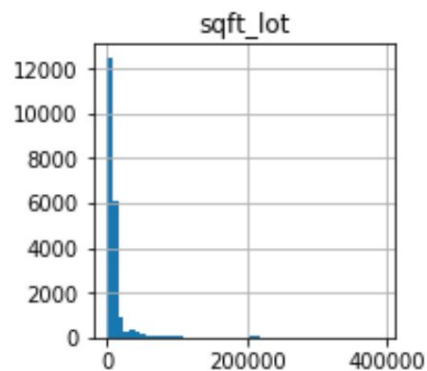
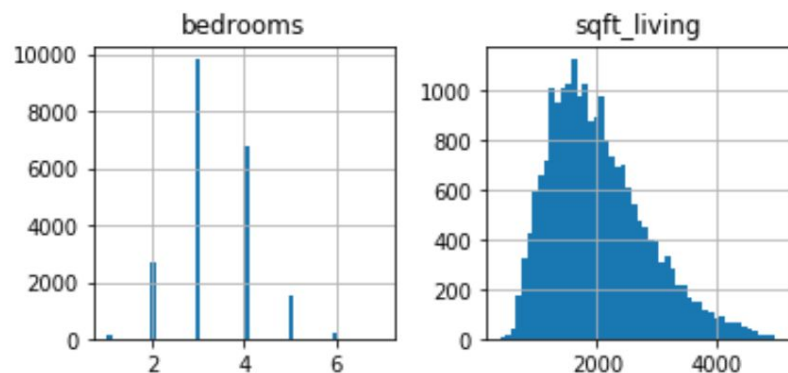
```
In [29]: df_dropped['sqft_lot'].corr(df['sqft_lot15'])
```

```
Out[29]: 0.7881861189991689
```

.75 Threshold



# Transforming a few features



Price  $\rightarrow$  Log of Price

