# MODELING SEATTLE HOUSING PRICES

by Matthew E. Parker

**Starting dataset:**

> 21,500 house sales
19 variables per sale

variables list:

```
price
bedrooms
bathrooms
sqft_living
sqft_lot
floors
waterfront
view
condition
grade
sqft_above
sqft_basement
yr_built
yr_renovated
zipcode
lat
long
sqft_living15
sqft_lot15
```

INPUTS

**Our approach to constructing a model:**
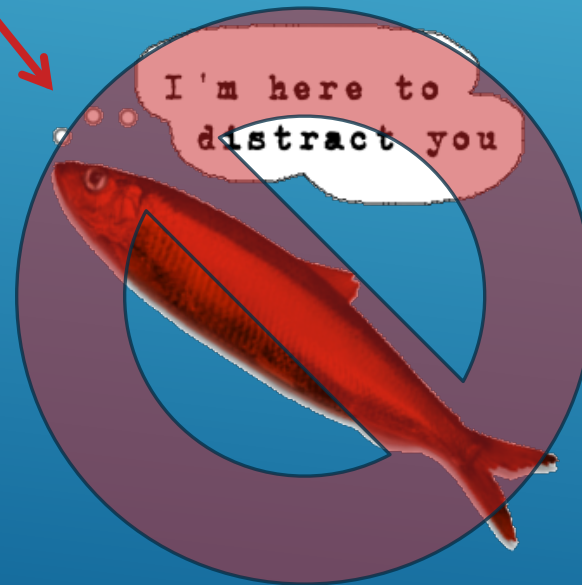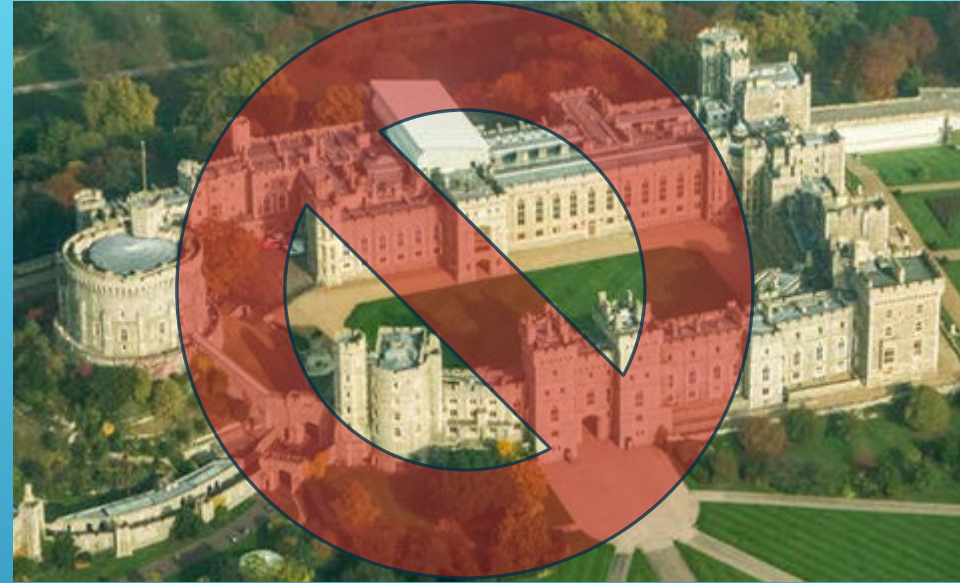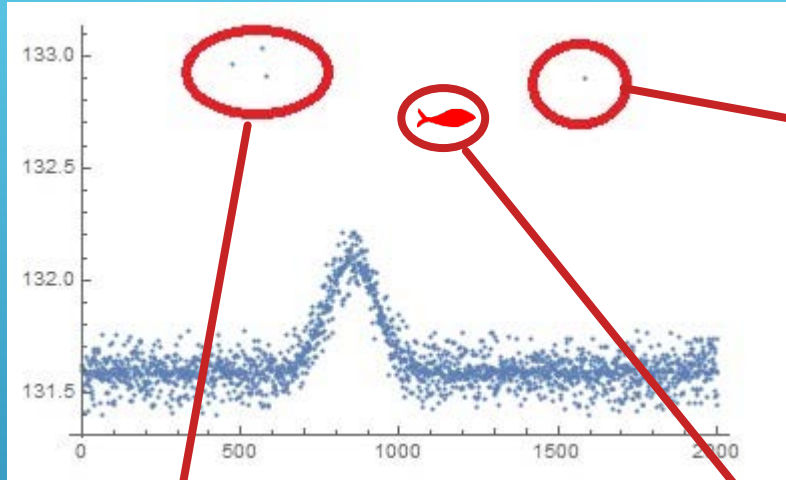
➢ Clean the data

➢ Explore and analyze the data

➢ Identify significant variables and build model around them

➢ Test and Validate model accuracy

Removing outliers from the dataset

CLEANING THE DATA
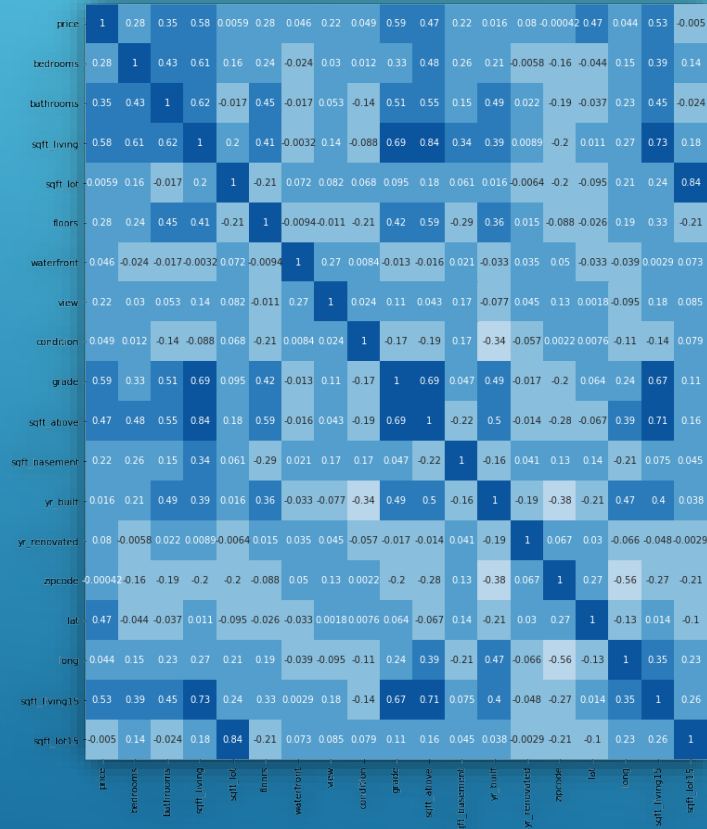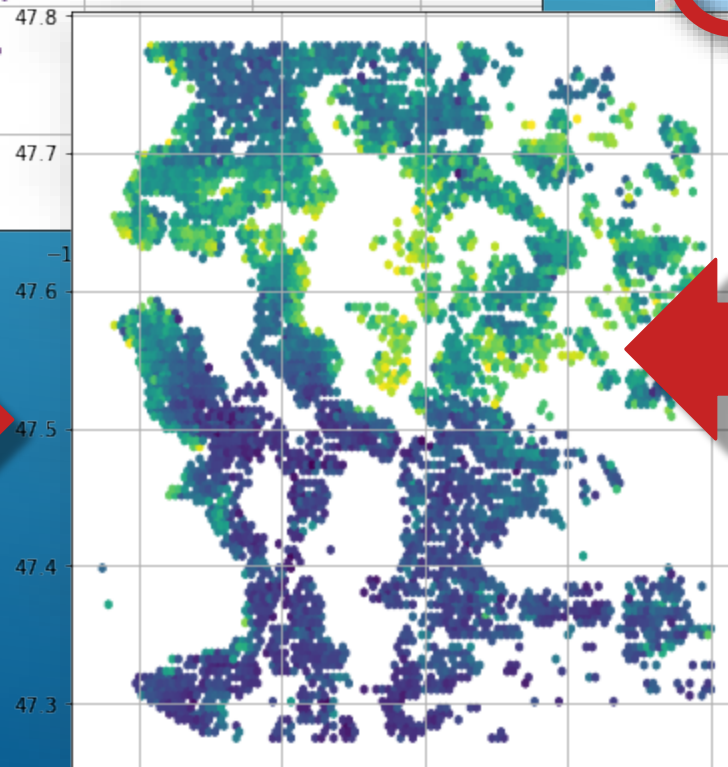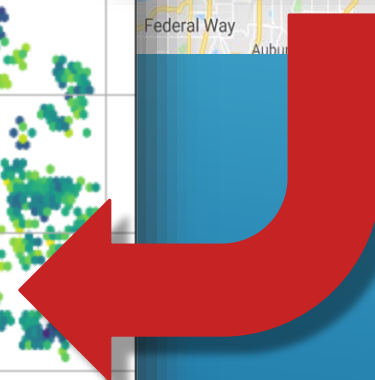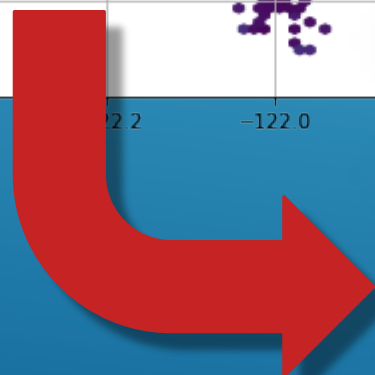
EXPLORING THE DATA
MAKING CONNECTIONS

Example of correlation matrix

Variables — Price

Using statistics to identify the variables with the greatest influence upon housing prices

HANDLING GEOGRAPHY

Total sq ft + Above sq ft + Below sq ft ≠ PRICE of a single HOUSE

Many variables are closely related, like **total sq ft**, **above-ground sq ft**, and **below-ground sq ft**.

Changing one variable can often impact another. This is bad for modeling as it produces a multiplied effect.

POTENTIAL PROBLEMS

**70%** Total sq ft $+$ **20%** Above sq ft $+$ **10%** Below sq ft $=$ **PRICE of a single HOUSE**

We don't want to just remove variables, as the presence or absence of a basement may influence pricing.
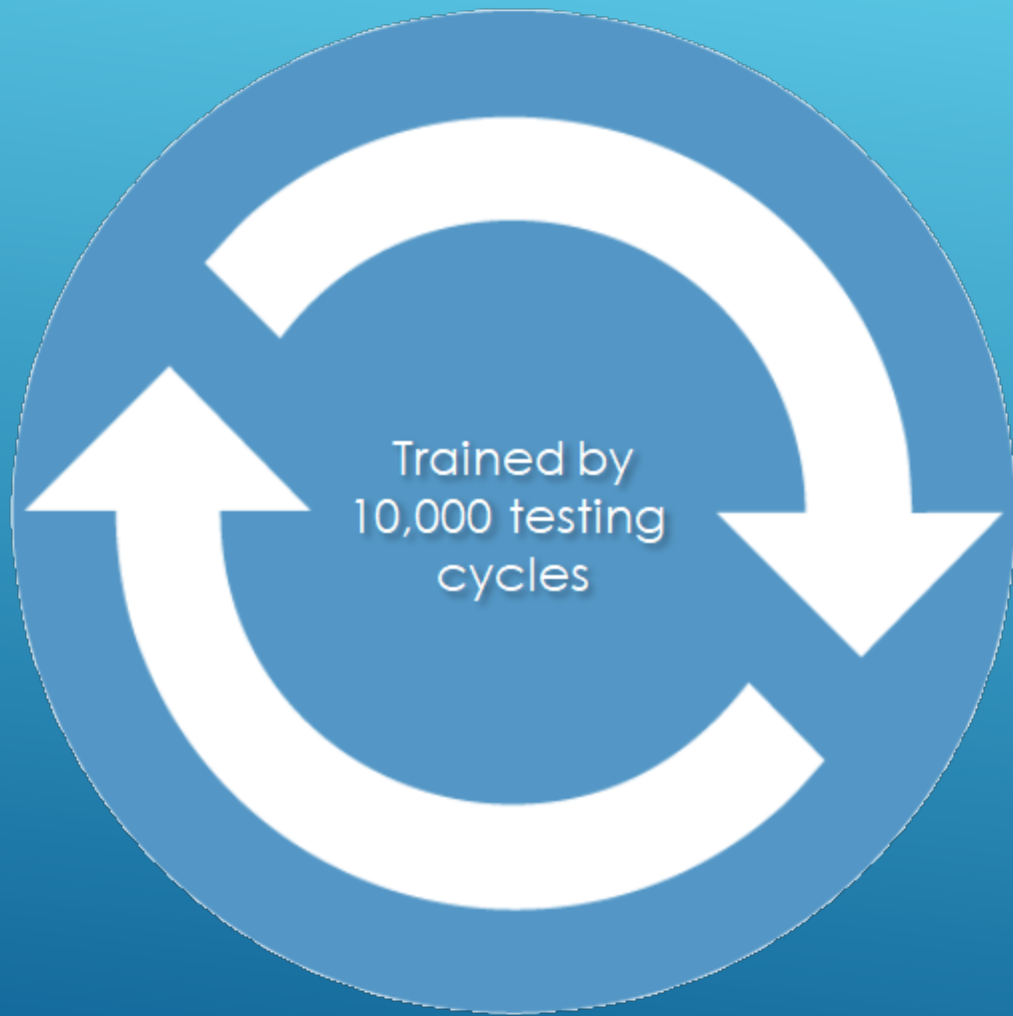
To solve this, we can build features that weight the variables proportionally to their influence.

BUILDING CUSTOM FEATURES

# TRAINING & VALIDATING THE MODEL



Trained by 10,000 testing cycles



Our model

$$p = (1296579.6 \times f_L) + (60401.32 \times f_V) + (110.79 \times f_B) + (59144.57 \times f_G) + (111362.3 \times \log_e(f_S)) - 62564620.4$$

$p$ = House price (in USD)

$f_L$ = latitude

$f_V$ = times property has been viewed

$f_B$ = square footage of basement

$f_G$ = grade given to the housing unit, based on King County grading system

$f_S$ = square footage of living space

## MODEL SUMMARY

If you know a house's latitude, basement ft$^2$, living space ft$^2$, King County grade, and the number of times it has been viewed, then you can estimate it's sale price within an error margin of $126,700.$^{00}$.

**Price =
Latitude +
Views +
Basement sq ft +
King County grade +
Total sq ft**

# RECOMMENDATIONS

1. If you can purchase a house for $126,700 less than the price predicted by our model, you will definitely make a profit.

2. Houses in the northern half of Kings County fetch higher prices, try to sell northern properties

3. The more times a house has been viewed, the higher it's final selling price is likely to be. Invest in advertising your properties.

# FURTHER INVESTIGATION

The current model could likely be enhanced by the addition of more variables.

In particular, information on crime rates, transportation accessibility, school district ratings, etc. would be useful as these factors have in the past been shown to influence real estate pricing.

# THANK YOU