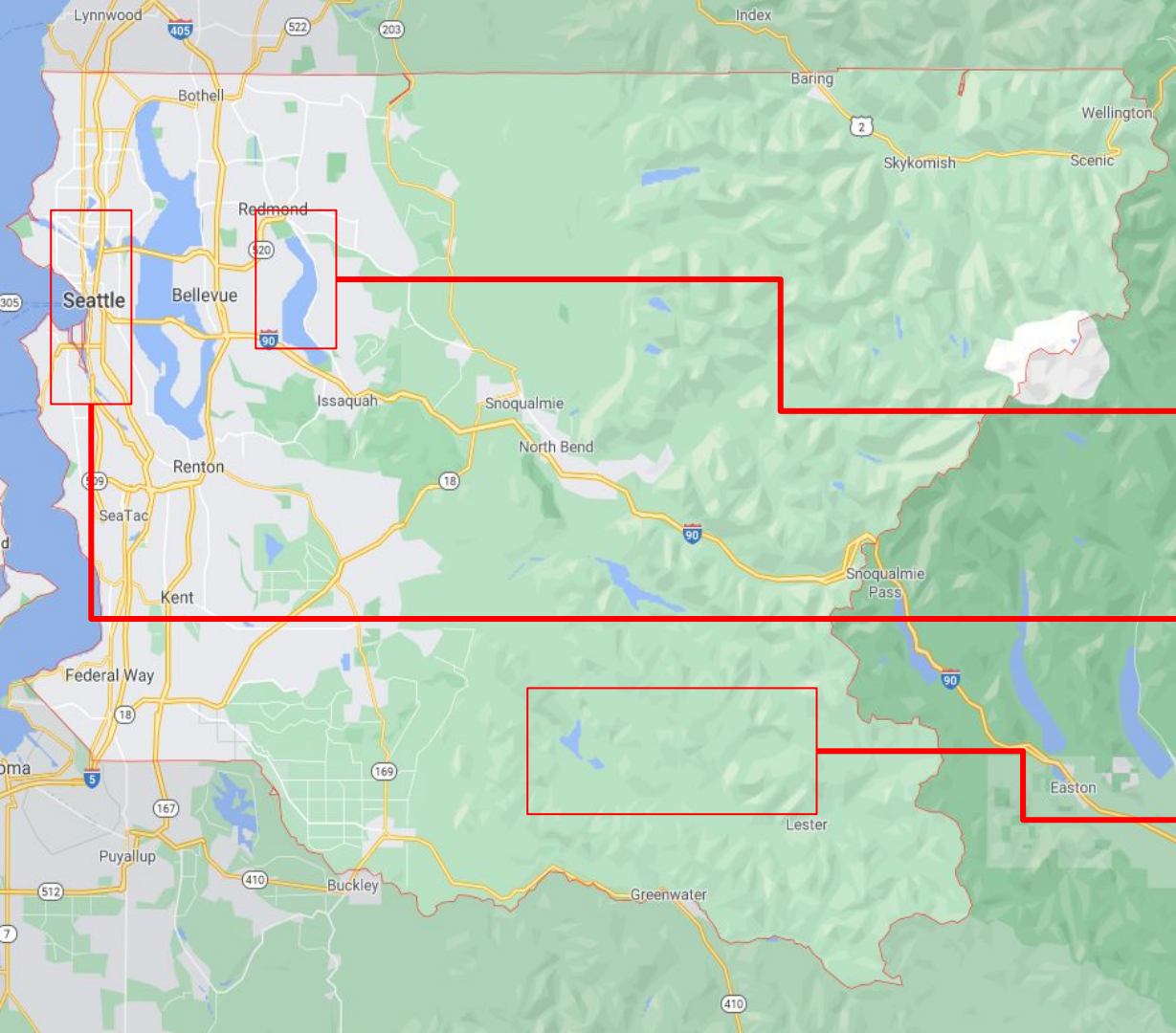# Predicting House Prices in King County, Washington

Leonard Boes

First look on map
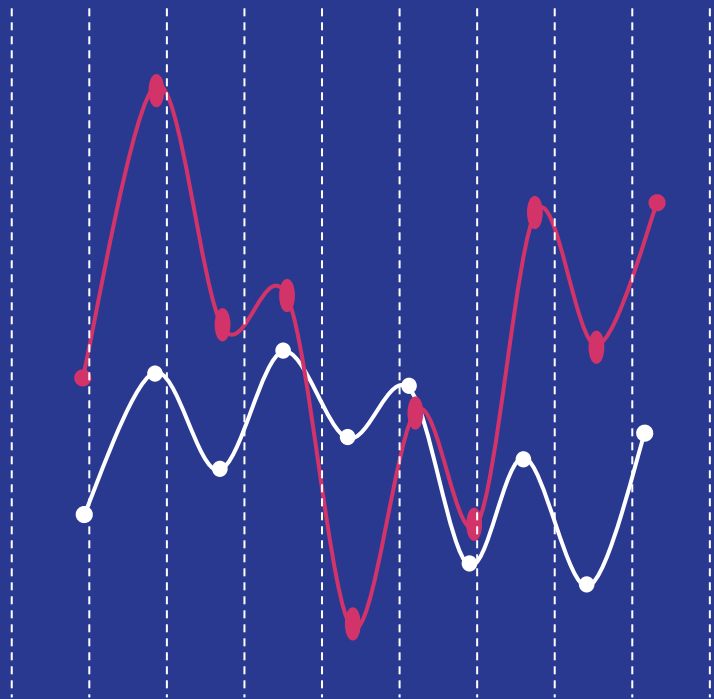
Unique areas like waterview

Dense Seattle Metropolis

Rural areas far away from city and interstate

2

# First look on the data

- Price and date of individual sales

- Number of bedrooms, bathrooms and floors

- Rank of condition and information about waterview

- Geographical location and zipcode

- Information about size of lot and living area

- Information about size of lot and living area of 15 neighboring houses

- Year of construction and renovation

- Roughly 21,600 observations

# Condition of the data

```
Data columns (total 21 columns):
id               21597 non-null int64
date             21597 non-null object
price            21597 non-null float64
bedrooms         21597 non-null int64
bathrooms        21597 non-null float64
sqft_living      21597 non-null int64
sqft_lot         21597 non-null int64
floors           21597 non-null float64
waterfront       19221 non-null float64
view             21534 non-null float64
condition        21597 non-null int64
grade            21597 non-null int64
sqft_above       21597 non-null int64
sqft_basement    21597 non-null object
yr_built         21597 non-null int64
yr_renovated     17755 non-null float64
zipcode          21597 non-null int64
lat              21597 non-null float64
long             21597 non-null float64
sqft_living15    21597 non-null int64
sqft_lot15       21597 non-null int64
```
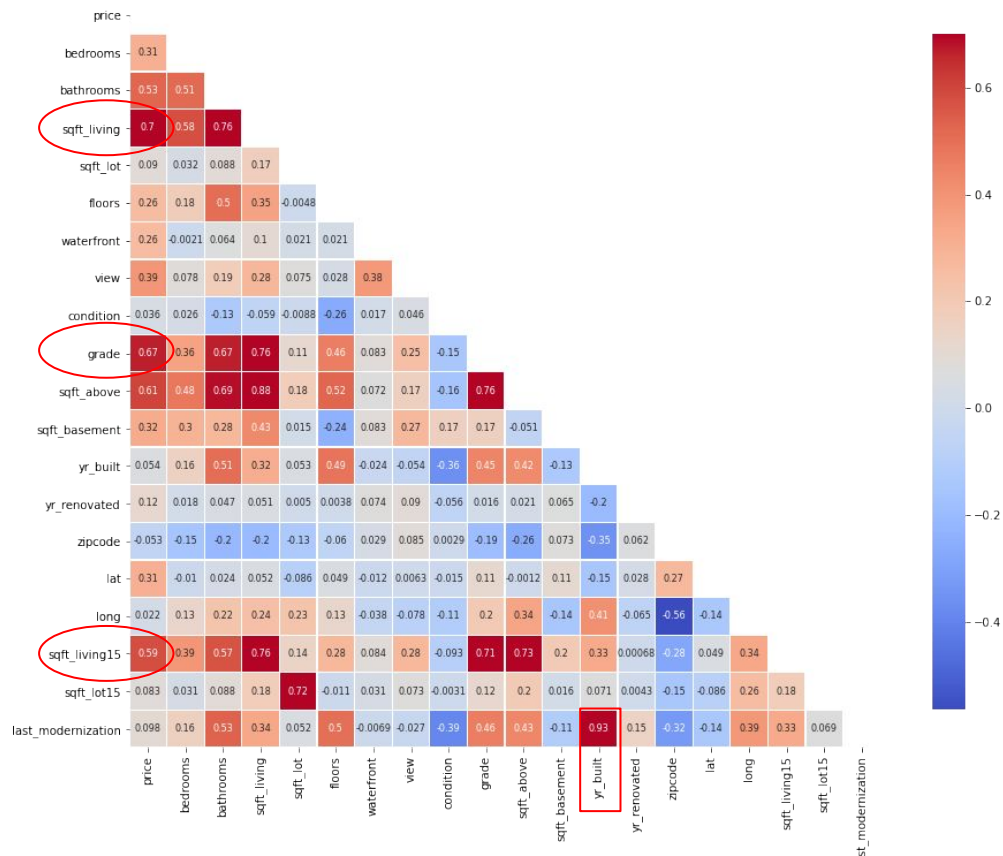
Date needed to be converted

Missing values needed to be accounted for

Something is wrong with the basement measures

Fortunately, only few values were missing and could logically be replaced with zeros
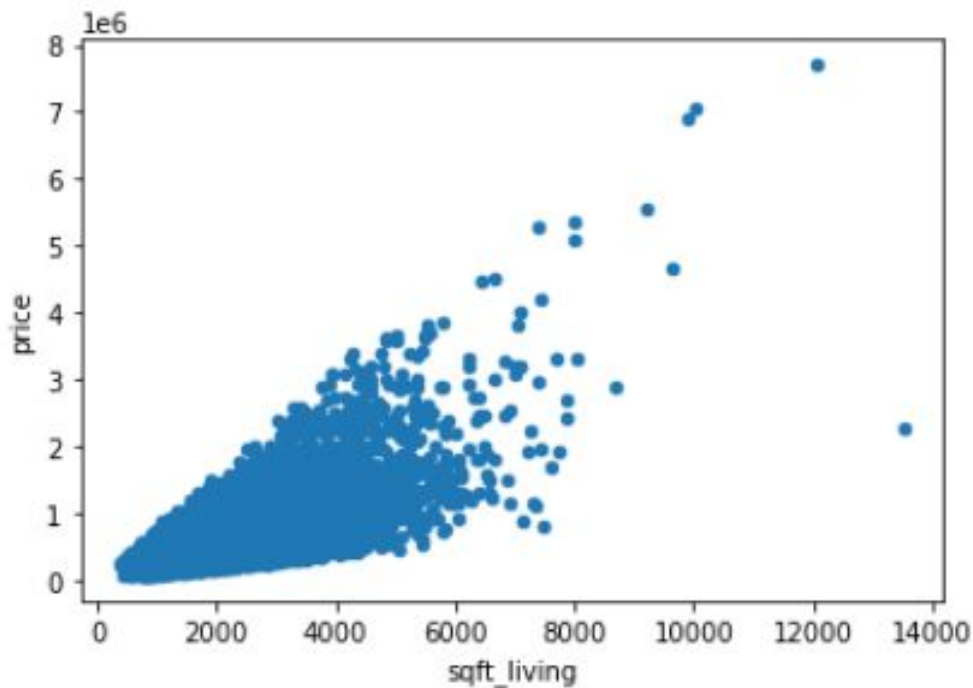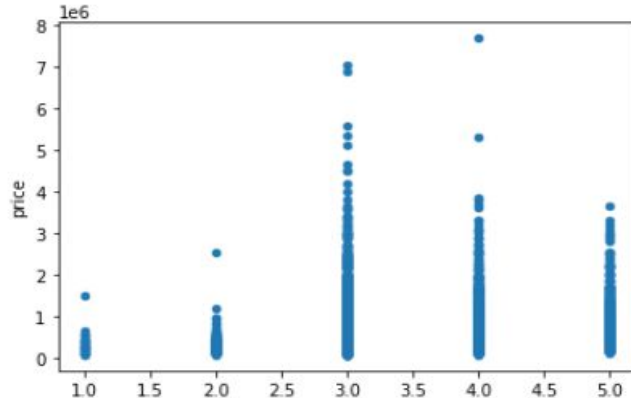
4

# Exploring the data

# Exploring the data



Highest correlation with price can be observed at:

Size of living area,
grade
& neighbors size of living area
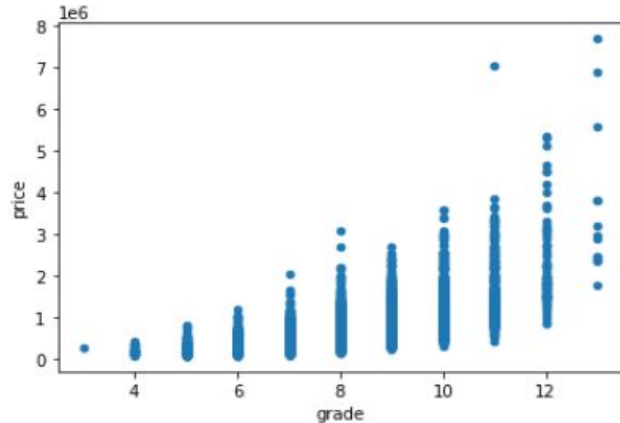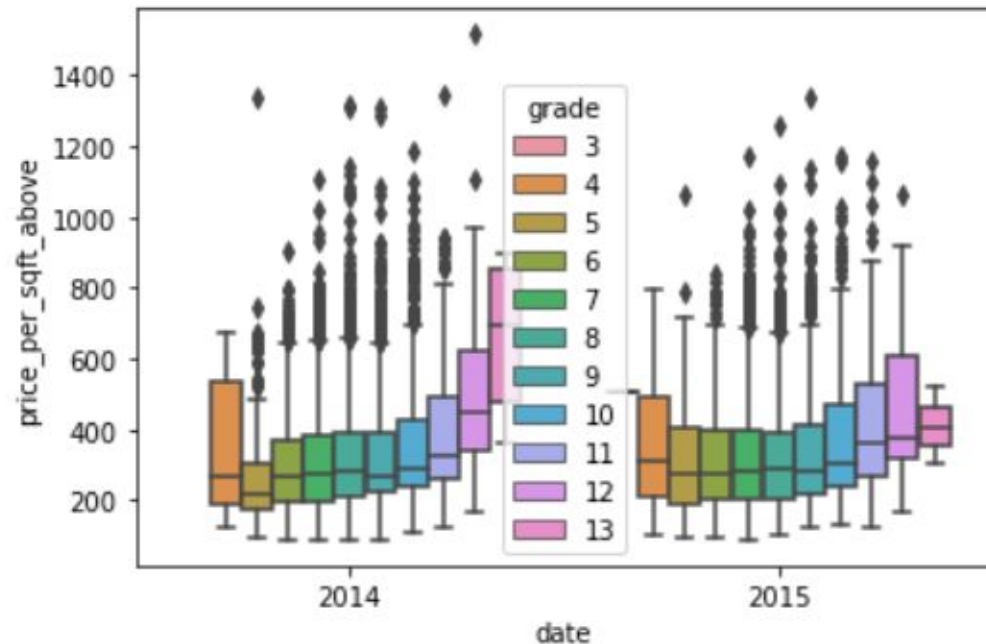
# Exploring the data



Highest correlation with price can be observed at:

Size of living area,
grade
& neighbors size of living area

# Exploring the data



No linear relationship between price and condition rating. Condition just needs top be > 2

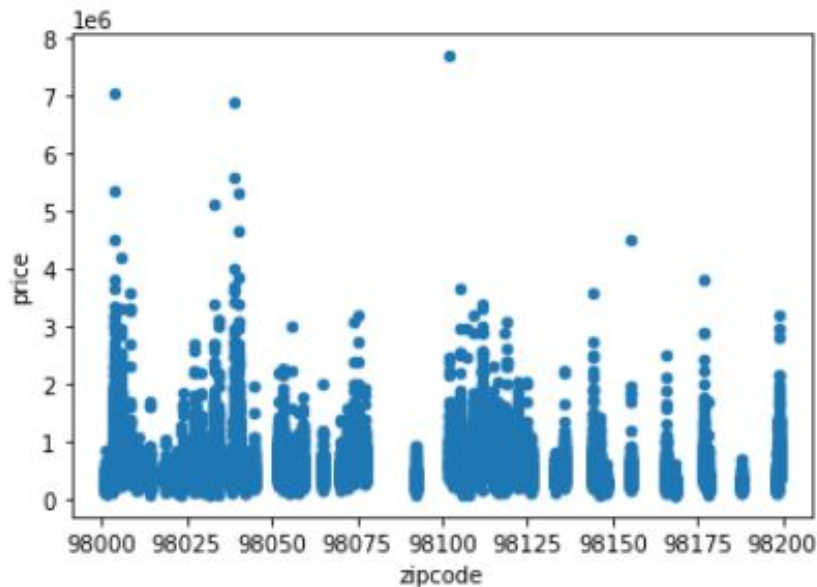Positive correlation trend between grade and price

# Exploring the data



If you want to buy luxury real estate, you might want to put a special emphasis on market timing

# Exploring the data



Some zipcode areas seem to have really cheap house prices.

The zip code areas with best price to footage ratio are:
98092, 98002, 98030, 98001, 98023, 98042, 98003 & 98038,
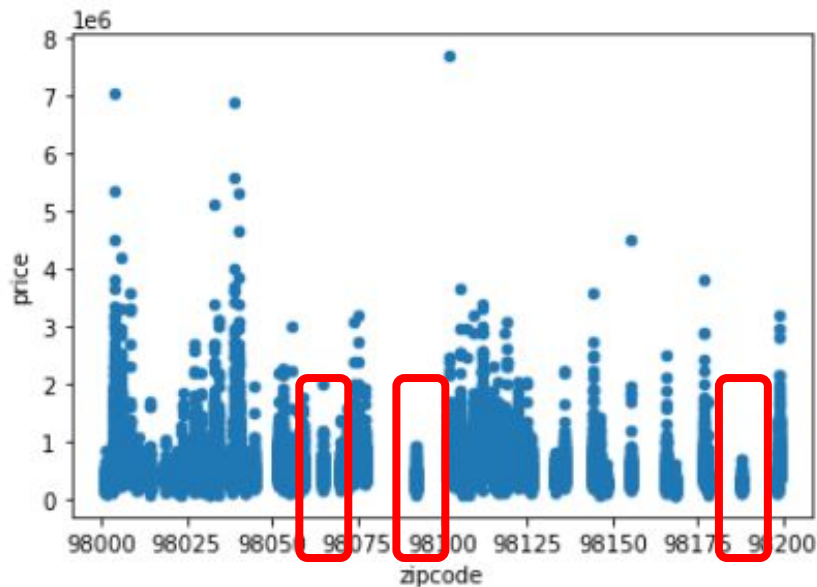
# Exploring the data



Some zipcode areas seem to have really cheap house prices.

The zip code areas with best price to footage ratio are:
98092, 98002, 98030, 98001, 98023, 98042, 98003 & 98038,

# Building new features

# Building new features

- Size of living area seemed to be important

- I wanted to add geographical dimension as well

```python
#calculate how many squarefeet the rooms have on average
df['rooms_per_sqft'] = (df.bedrooms + df.bathrooms) / df.s

#calculate how many bedrooms per squarefoot are there
df['bedrooms_per_sqft'] = df.bedrooms / df.sqft_living

#calculate how many bathroomy per squarefoot there are
df['bathrooms_per_sqft'] = df.bathrooms /df.sqft_living

#create lot to above ratio for individual home
df['lot_above_ratio'] = df.sqft_lot / df.sqft_above

#create lot to living ratio for individual home
df['lot_living_ratio'] = df.sqft_lot / df.sqft_living

#create lot to living ratio for nearest 15 neighbors
df['lot_living_ratio15'] = df.sqft_lot15 / df.sqft_living1

#calculate how many rooms you get for a dollar
df['bedrooms_per_dollar'] = df.bedrooms / df.price
```
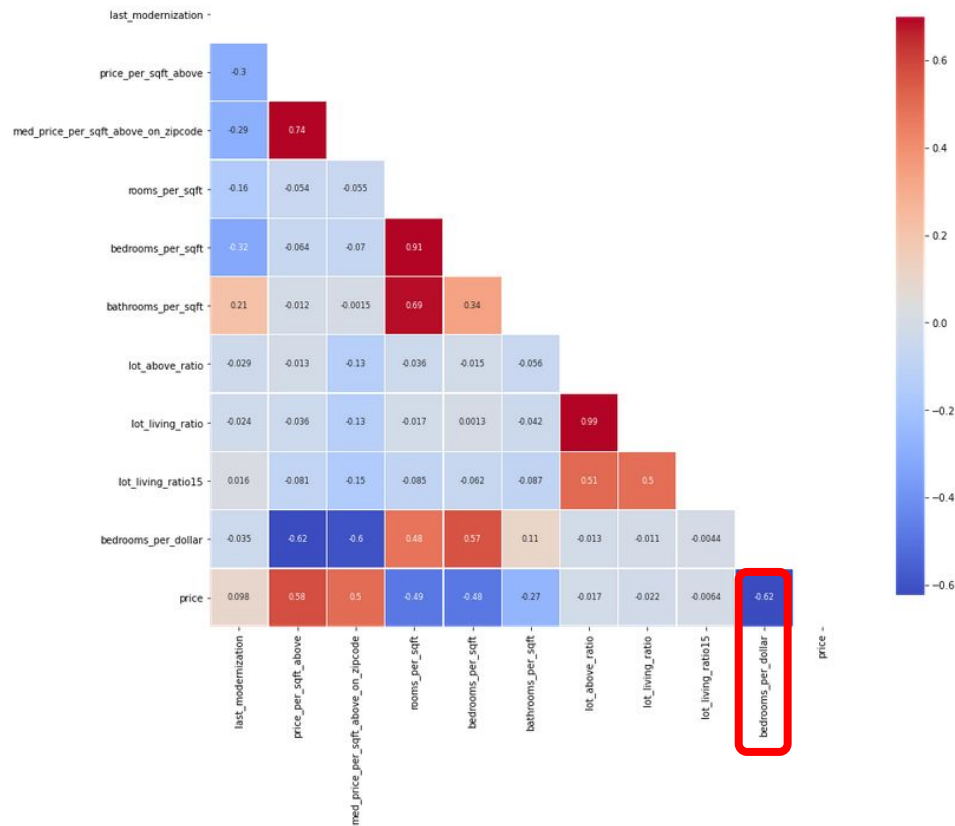
# Exploring the data



The more bedrooms you get for a dollar, the cheaper the house price.

Sellers should keep the number of rooms normal and not try to squeeze more rooms in a property that neccessary

# Building a predictive model

# Looking at different model qualities

## Best Single Predictor

**Square footage of living area**

RSquared: 0.469
RMSE: $ 249,948

Price = -52,451 + 285 x area

## Best Feature Combination

**Square footage of living area, grade & year of construction**

RSquared: 0.598
RMSE: $ 217,556

Price = 6319009
        + 182 x area
        + 144080 x grade
        -  3684 x year of construction

# Looking at different model qualities

**Old Best Feature Combination**

**New Best Feature Combination**

**Square footage of living area, grade & year of construction**

**Square footage of living area, median price per sqft in zip-area, waterfront**

RSquared: 0.598
RMSE: $ 217,556

RSquared: 0.758
RMSE: $ 168,863

Price = 6319009
+ 182 x area
+ 144080 x grade
- 3684 x year of construction

Price = -497856
+ 267 x area
+ 1594 x median price per sqft in zip area
+ 839,473 x waterfront

17

# Looking at different model qualities

**Old Best Feature Combination**

**Square footage of living area, grade & year of construction**

RSquared: 0.598
RMSE: $ 217,556

Price = 6319009
        + 182 x area
        + 144080 x grade
        - 3684 x year of construction

**New Best Feature Combination**

**Square footage of living area, median price per sqft in zip-area, waterfront**

RSquared: 0.758
RMSE: $ 168,863

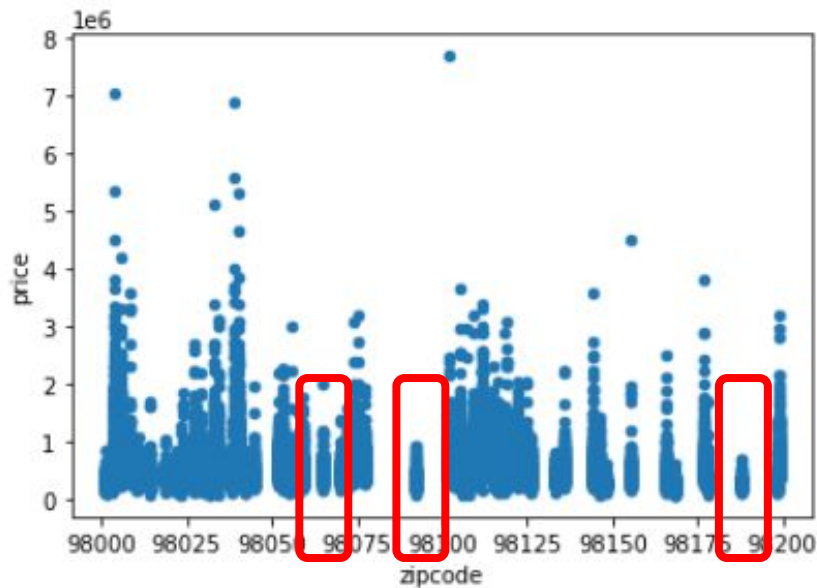Not the margin of error a house buyer wants to see

Price = -497856
        + 267 x area
        + 1594 x median price per sqft in zip area
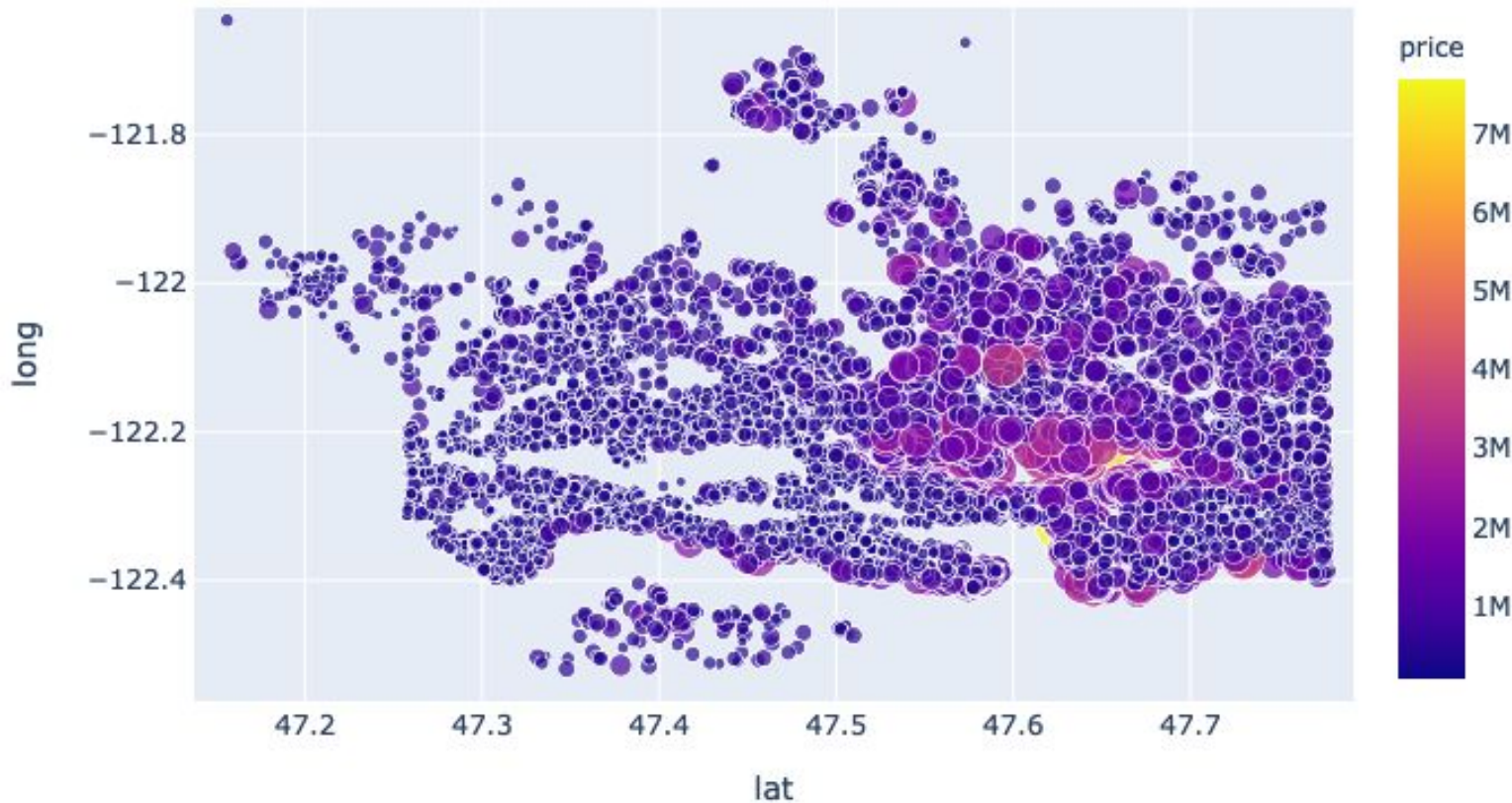        + 839,473 x waterfront
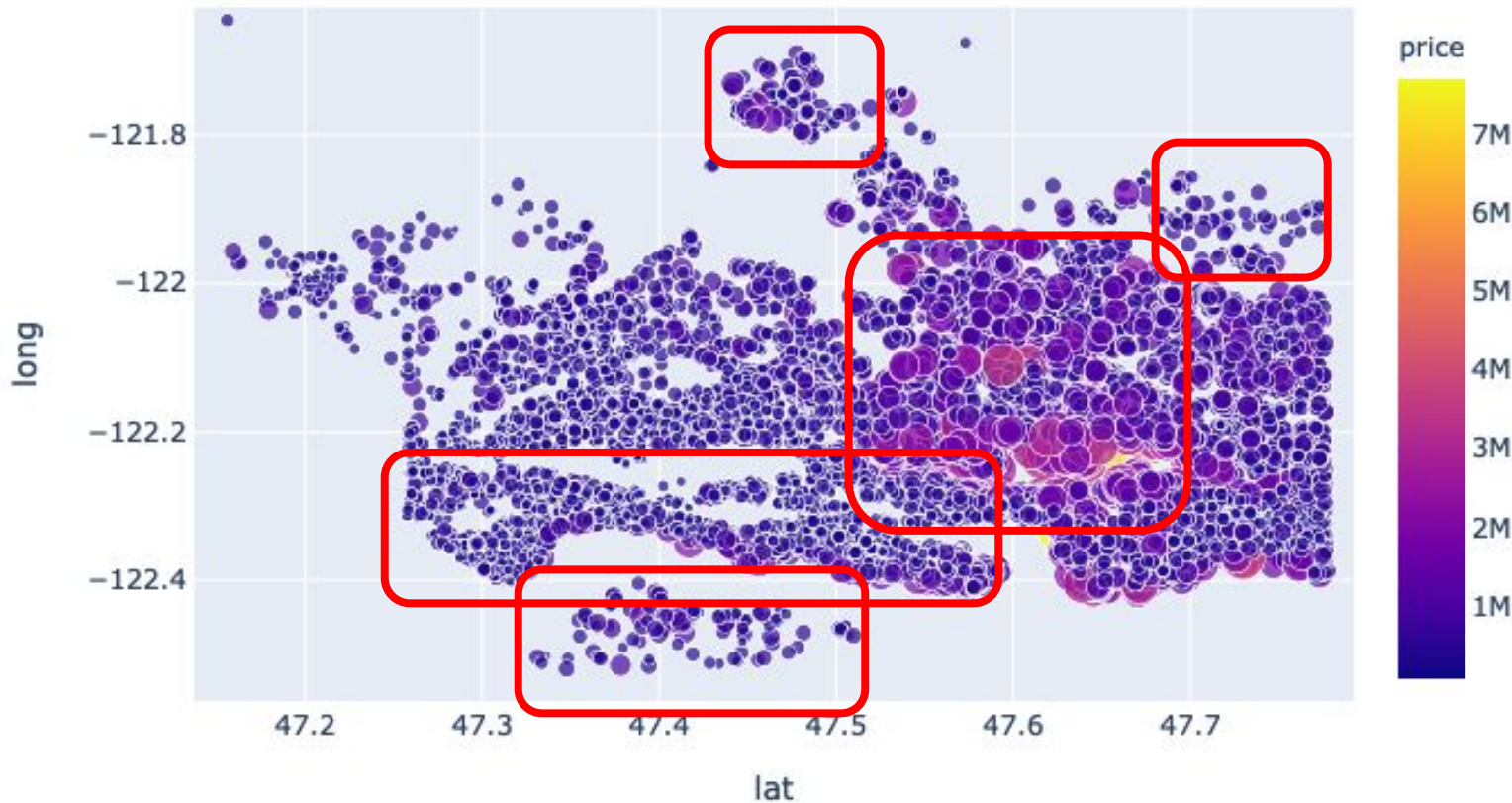
# Outlook

# Remember this chart?



Some zipcode areas seem to have really cheap house prices.

# Clustering by region could be even better feature

# Clustering by region could be even better feature

# Thanks for your attention