# 92173

A Data Science Perspective
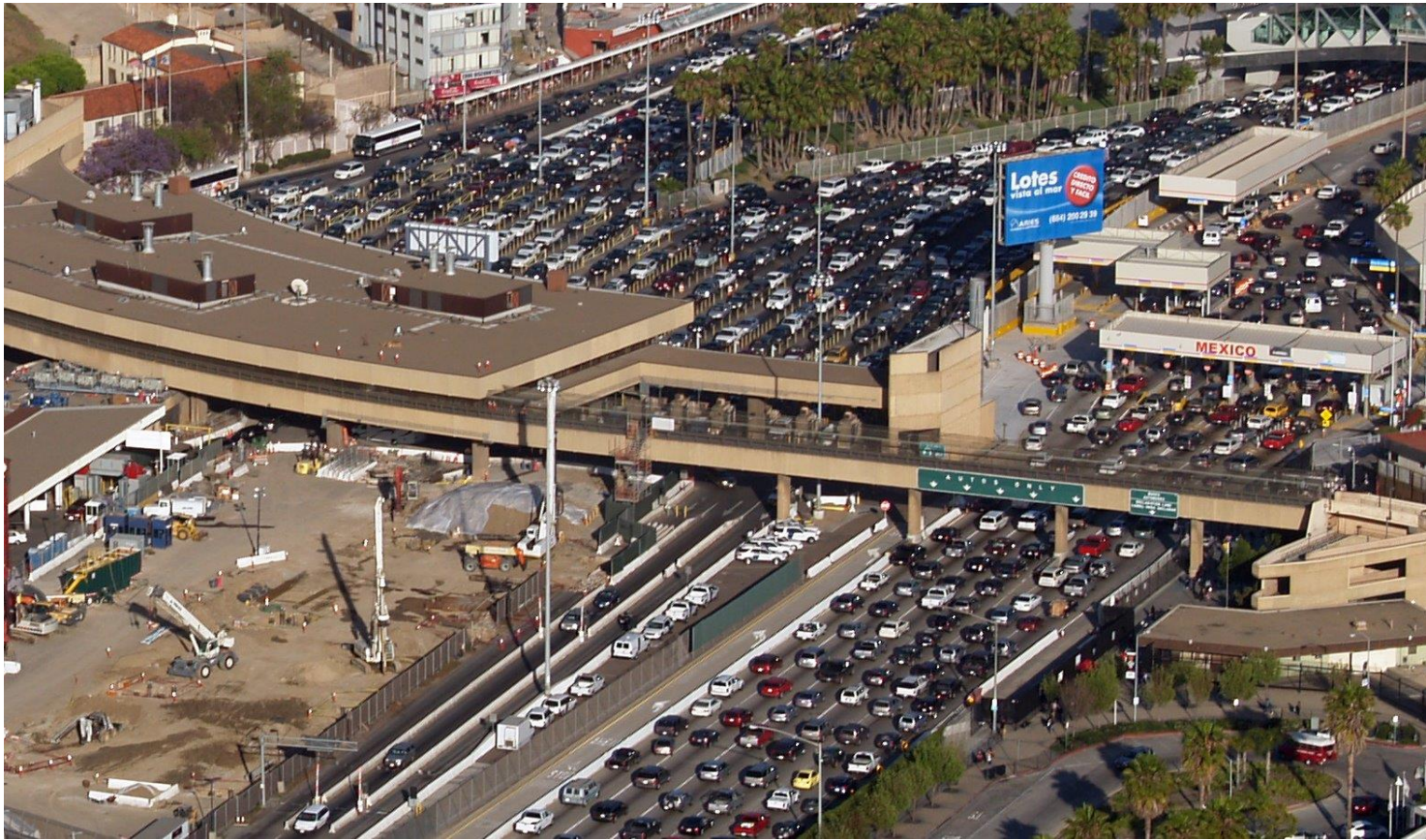
# Scope: Visualization and Machine Learning

We will start by doing a real estate analysis of the 92173 Zip code.

This is a Zip code very dear to me as I have lived here for more than 10 years.

Code will be provided on a GitHub link.

- The city of San Ysidro has been around since the 1900's.

- Serves as a gate between Mexico and the United States.

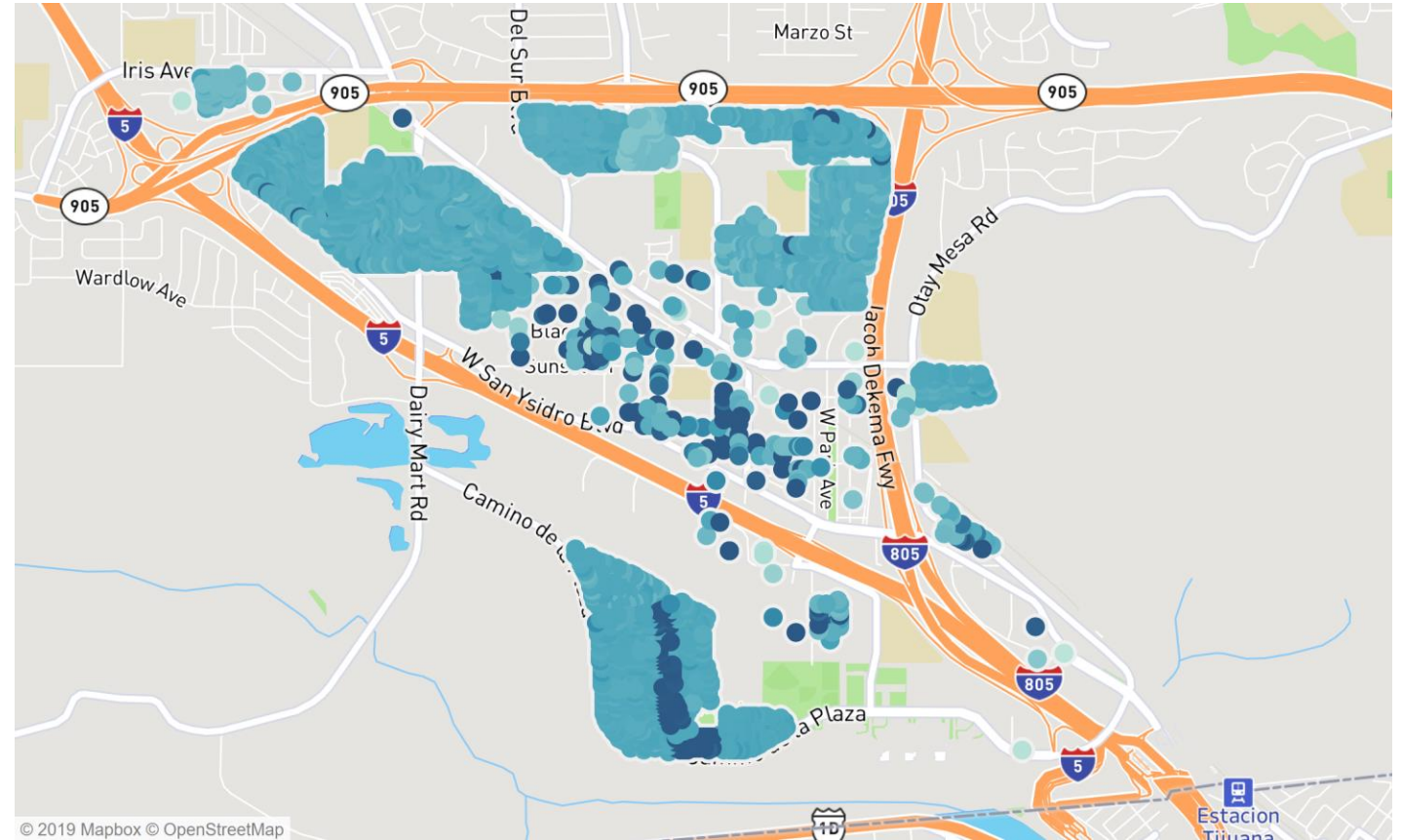- Also a rich community with a 93.8% Hispanic demographic makeup.

# Between Two Worlds

Source:https://statisticalatlas.com/neighborhood/California/San-Diego/San-Ysidro/Race-and-Ethnicity

# Initial Look at Housing

- After connecting to the Zillow Api I managed to scrap data for 2415 residencies.

- They are color coded by price, the darker the color the higher the price.

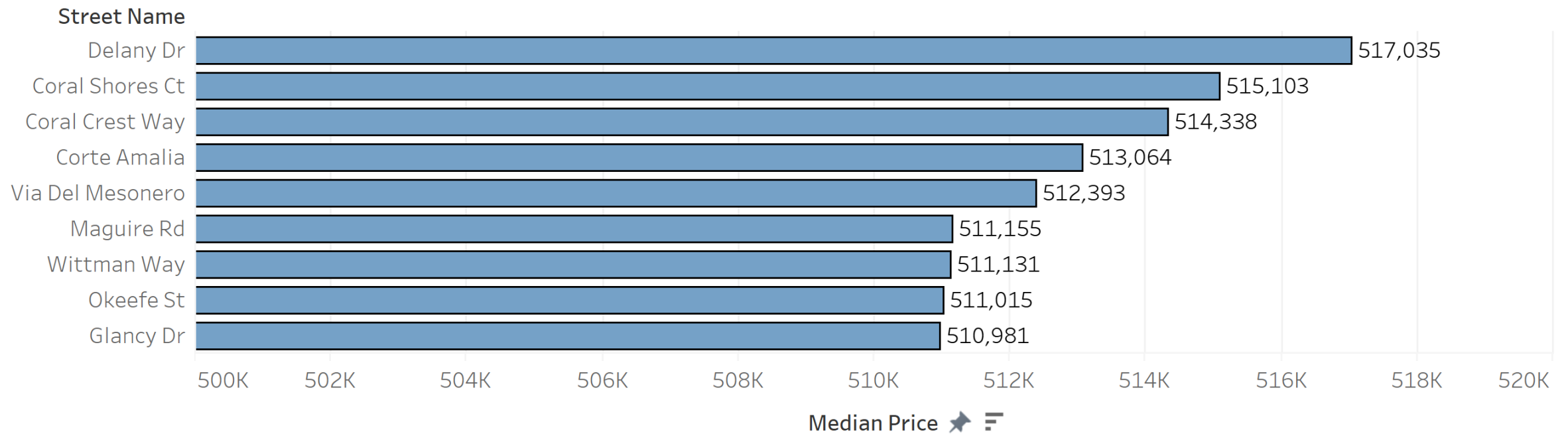- Full interactive map can be found here:
  - https://qrgo.page.link/tGB89

San Ysidro House Map by Prices
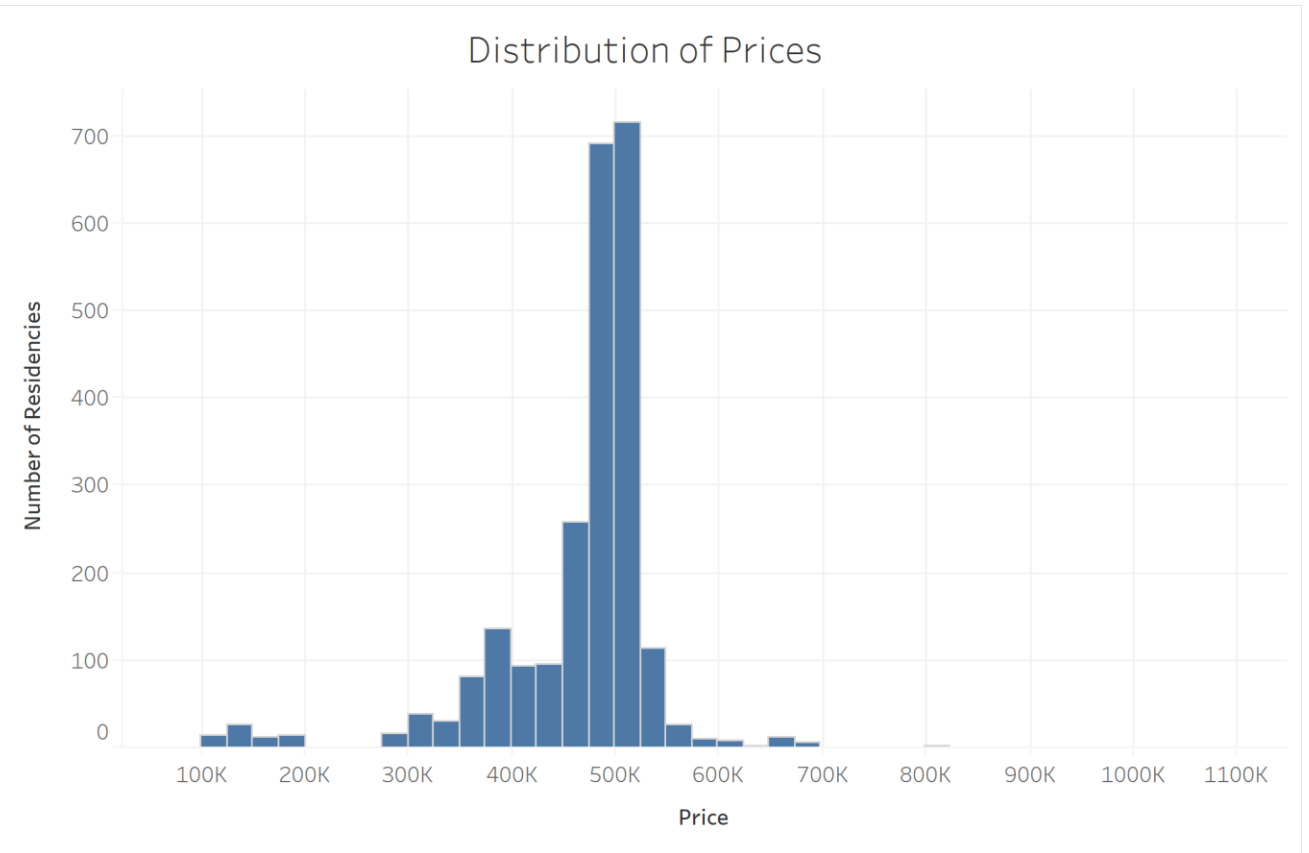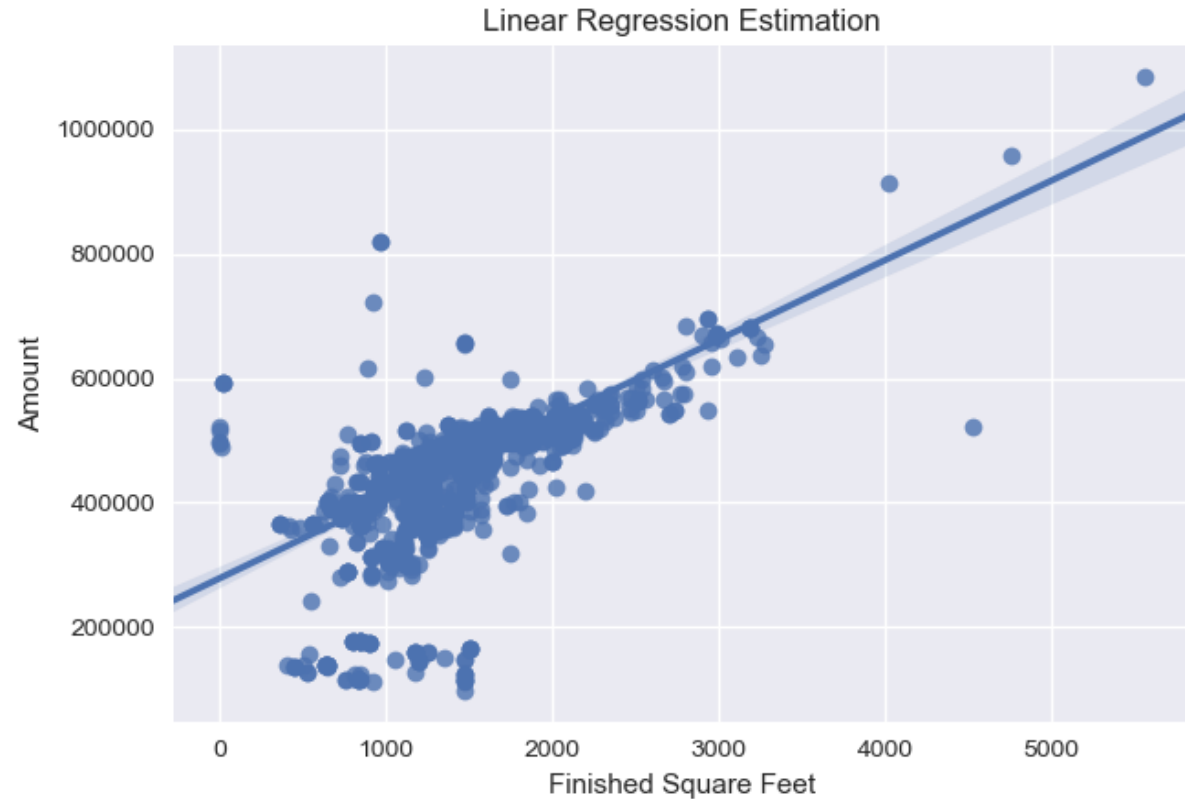
# Distribution

- Constructing a histogram we can take away the following.
    - Most houses are between 450K – 550K USD.
    - There is a small cluster of residencies from 100K to 200K USD.
    - This is follows a normal distribution.



Distribution of Prices
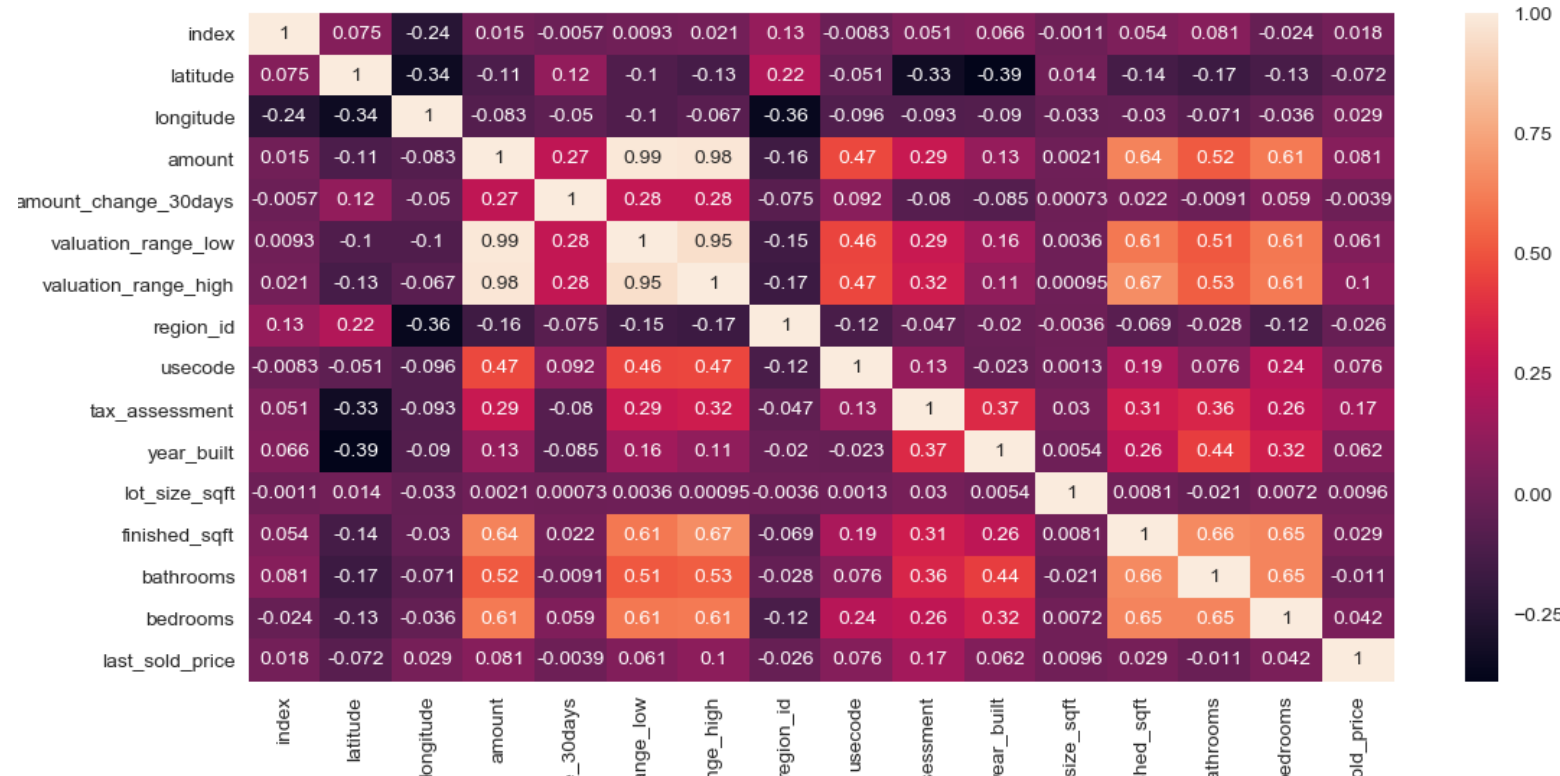
# Regression Models (Price)

- This type of continuous data on Prices requires regression techniques to make estimates.

- On the right we can see a method known as Linear Regression minimized by OLS.

- We draw the line of best fit and predict values based on that.



Linear Regression Estimation

# Feature Selection and Correlation

- Before we perform regression techniques for prediction, we have to pick which features are most important.

- A good step is to draw a correlation matrix and choose

- We chose:
  - Usecode
  - Finished_sqft
  - Bathrooms
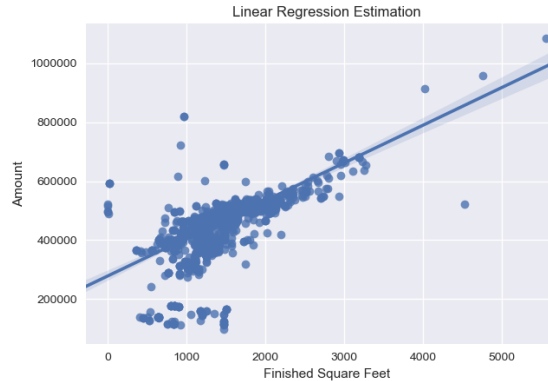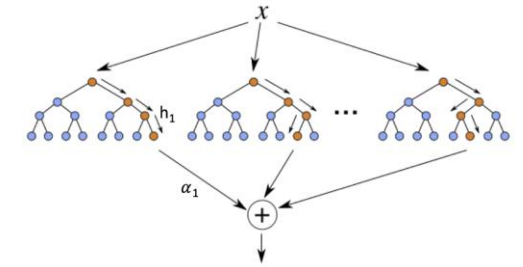  - Bedrooms

## Correlation Matrix

# Machine Learning Selection

- We want to now guess and be able to predict the log(price) of a home with our selected features.

- We Will use three different models and see how they each compare.

- Three models:
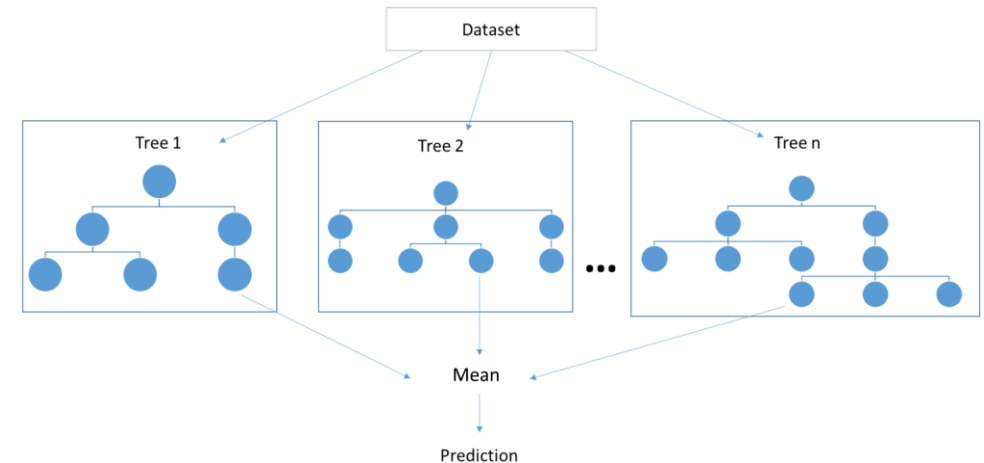  - Linear Regression
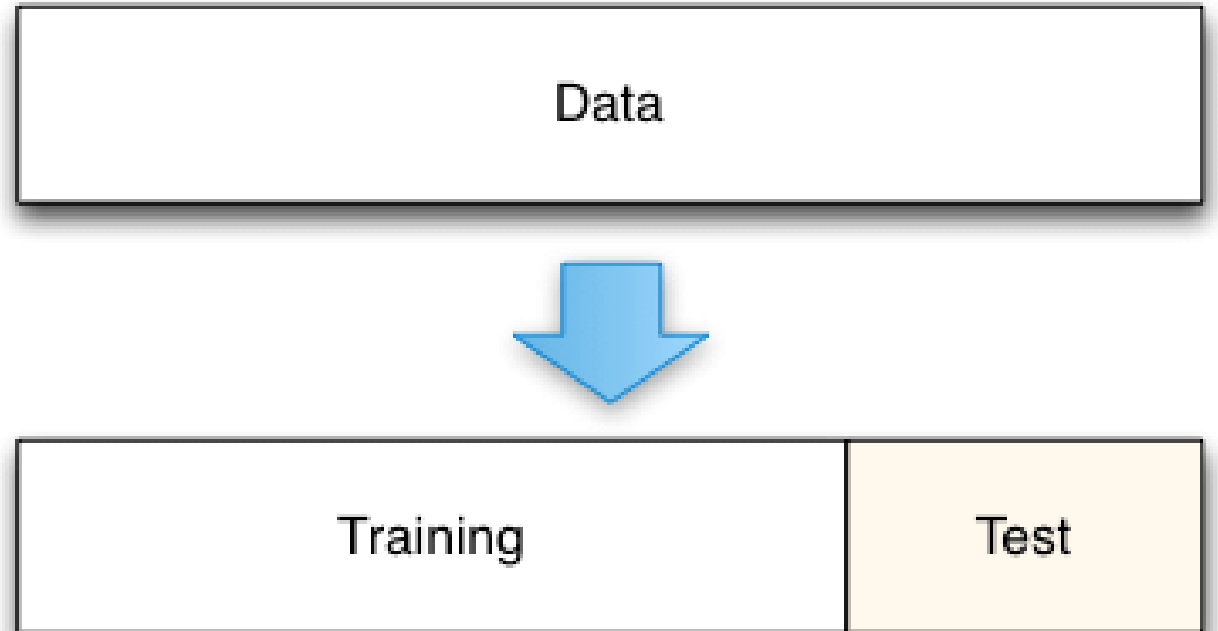  - Random Forest
  - XGBoost



Linear Regression



XGBoost



Random Forest

# Training and Testing

- We split our data into two data sets training and testing, 70% and 30% respectively.

- This is done to test our machine learning model on observations it theoretically has not seen before.

| Data |
|------|

⬇

| Training | Test |
|----------|------|

# Machine Learning Results

After our results we can see that the lowest absolute error comes from a random forest regression of 0.090. We will interpret this in the next slide

| Target Variable | Log(Amount) | | |
|---|---|---|---|
| Model | Linear Regression | **<u>Random Forest</u>** | XGBoost |
| MAE | 0.105 | **<u>0.090</u>** | 0.102 |

# Interpretation

Our best performing model was our Random Forest Regressor which scored a Mean Absolute Error of about 0.09 being the smallest error.

This mean that our model was approximately 9% off from the geometric mean.

This is an alright score and we theoretically should have scored higher ,but there are some issues within our data gathering process.

# Analysis Problems

To improve further models we can think about our data gathering. The Zillow API was constantly missing observations of mobile homes and apartment complexes.

These two types make up a great deal of the San Ysidro community and not including them raised the overall mean of the zip code.

There was a systematic bias in our data collection which in turn is reflected in our models.