

Predicting Home Values Through Random Forest

BY: ALFREDO MARTINEZ

2018

A solid orange horizontal bar spanning the width of the slide at the bottom.

Introduction

Goal

- Build a Predictive Model System for Home Prices

Client

- Keller Williams Realty, Inc.
- Other Real Estate Investment Companies



Exploratory Data Analysis(EDA)

Housing Dataset

Categorical(Qualitative) data:

- 23 nominal (There was no natural order, e.g. Type of house)
- 23 ordinal (Order do exist, e.g. Property condition; bad, fair, good excellent)

Numerical data:

- 14 discrete(Integers, e.g. Number of rooms)
- 20 continuous(Can take on any value, e.g. Square Feet)

Weather Dataset

- Subtracted from the National Center for Environmental Information Webpage.
- Includes daily, monthly and yearly records from different weather stations in the city.

Exploratory Data Analysis(EDA)

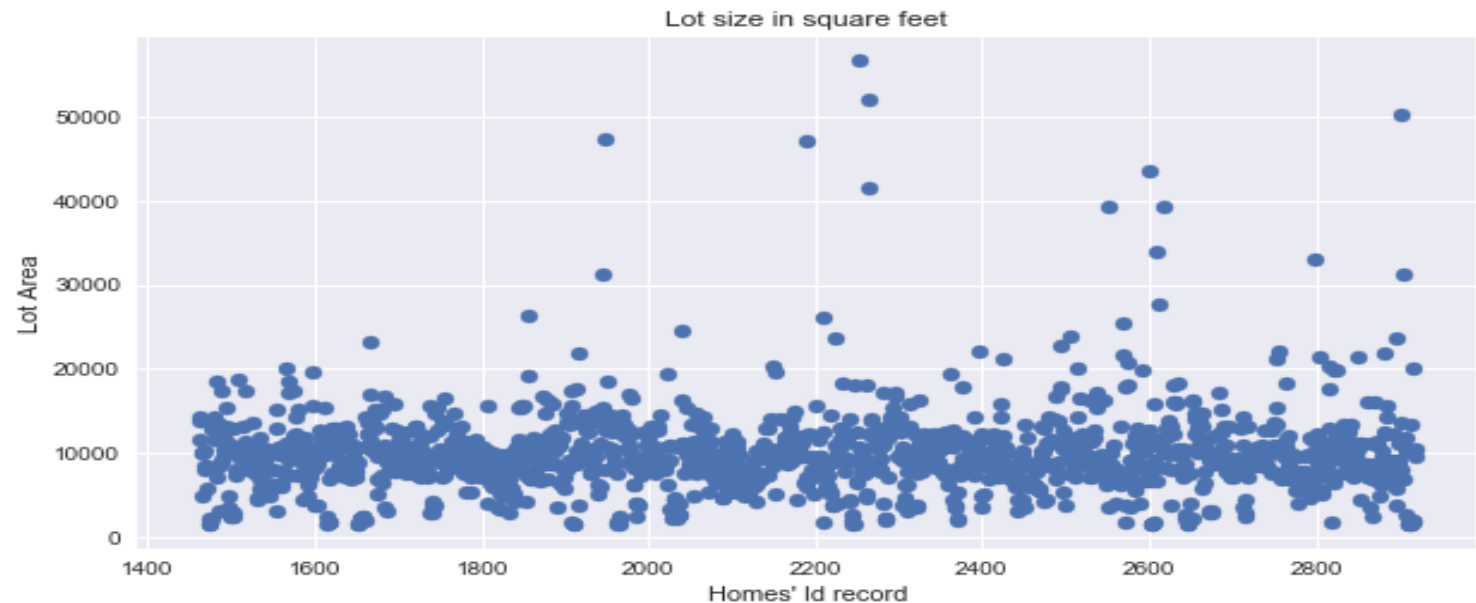
Bad Data and Outliers

Random Forest:

Different people, different arguments weather to keep all data or not

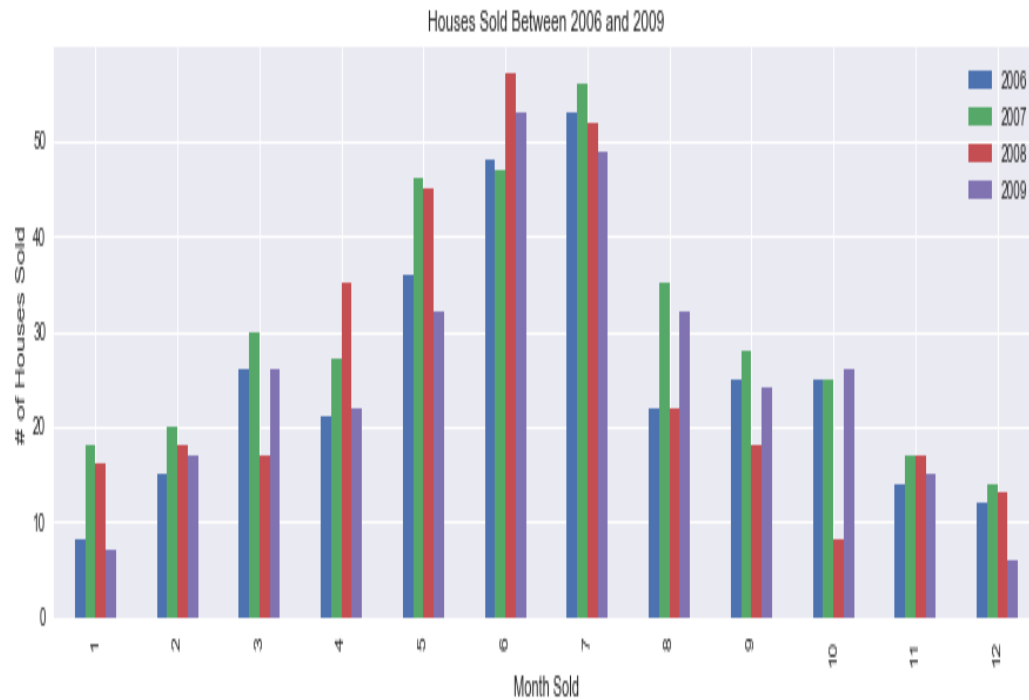
Ideas

- Not sensitive to outliers
- Sensitive



Exploratory Data Analysis

MONTHLY HOMES' SALES



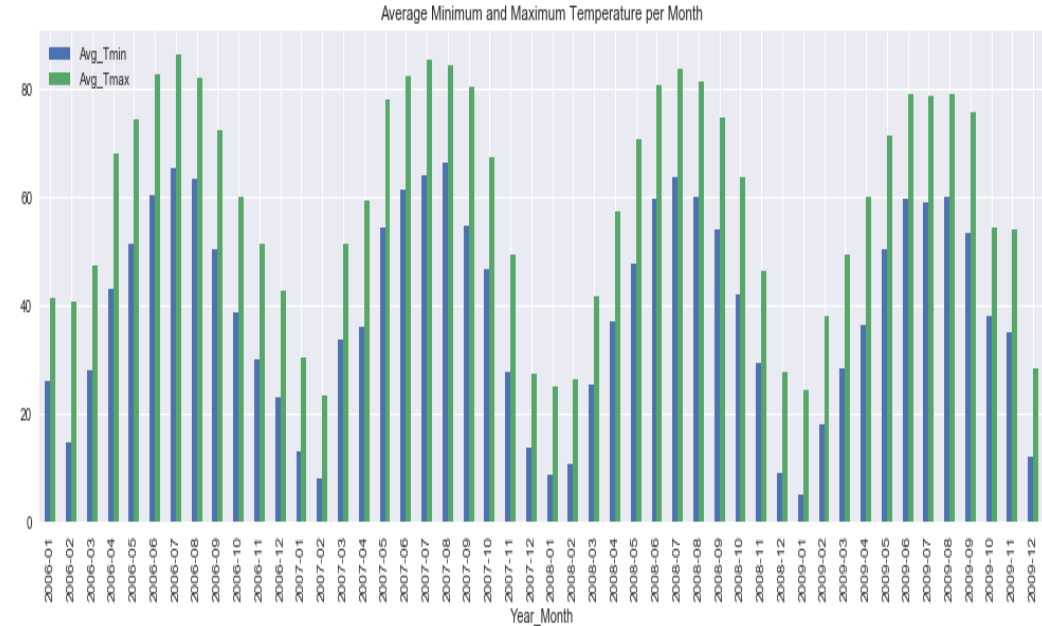
ABOVE GROUND LIVING AREA IN RELATIONSHIP TO SALE PRICE



Exploratory Data Analysis

HOME'S LOT AREA IN RELATIONSHIP TO PRICE

YEARLY WEATHER TEMPERATURE



Inferential Analysis

Alpha set at .05

Significant price difference analysis about population.

Central unit vs not a central unit

P-value: 0.010090621217186681 (**Significant**)

One story houses vs Two story houses

P-value: 3.5535259636604621e-13 (**Significant**)

Good privacy fence vs minimum privacy fence

P-value: 0.37003494954970051 (**Needs further research**)

Data Challenges

Missing Data

- N/A Values

Corrupted Data

- Present as (-999) values

Categorical to Numerical

- Transforming values to be able to use in predictive model(scikit-learn)

Fixed Values

- Mean, Median, Mode, and dummy variables indexing(from train to test)

Technology

Pandas

- To open and merge files
- To filter and drop cells
- To group by cells
- Etc..

Numpy

- To fill in Not a number values(NaN)
- Etc..

Datetime

- To fix and edit date times

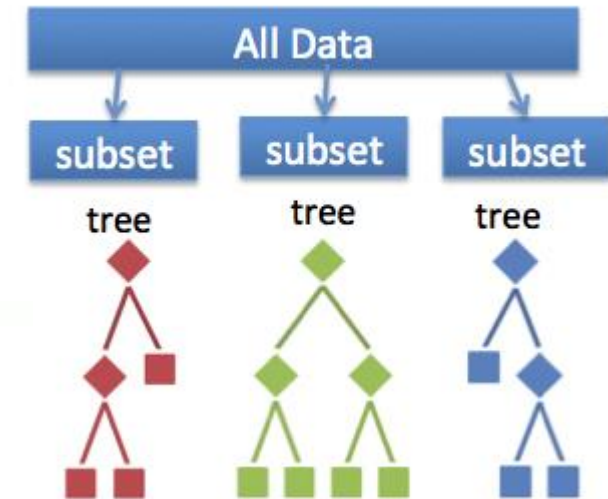
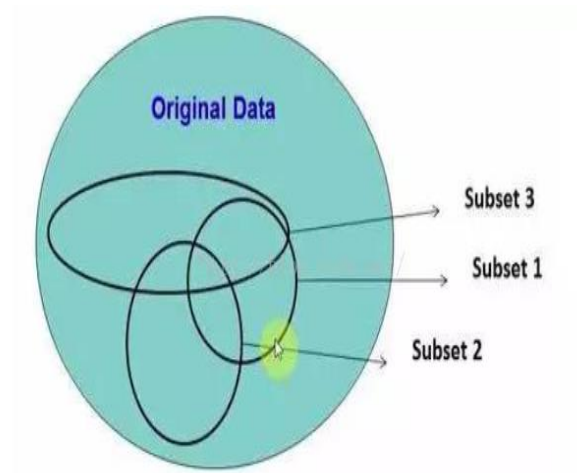
Scipy

- T-test for inferential analysis

Modeling and Analysis

Random Forest Predictive Algorithm

- Continues Target variable(Regression)
- Based on Decision Trees Models
- Bootstrapping
- Root Mean Square Error Accuracy
- Feature Importance



Modeling and Analysis

Main packages and libraries

Scikit-learn and math

- `RandomForestRegressor()`
- `Param_grid(GridSearchCV)`
- `Fit()`
- `Best_estimators()`
- `Predict()`
- `RMSE = sqrt(mean_squared_error(y_test, y_predict))`

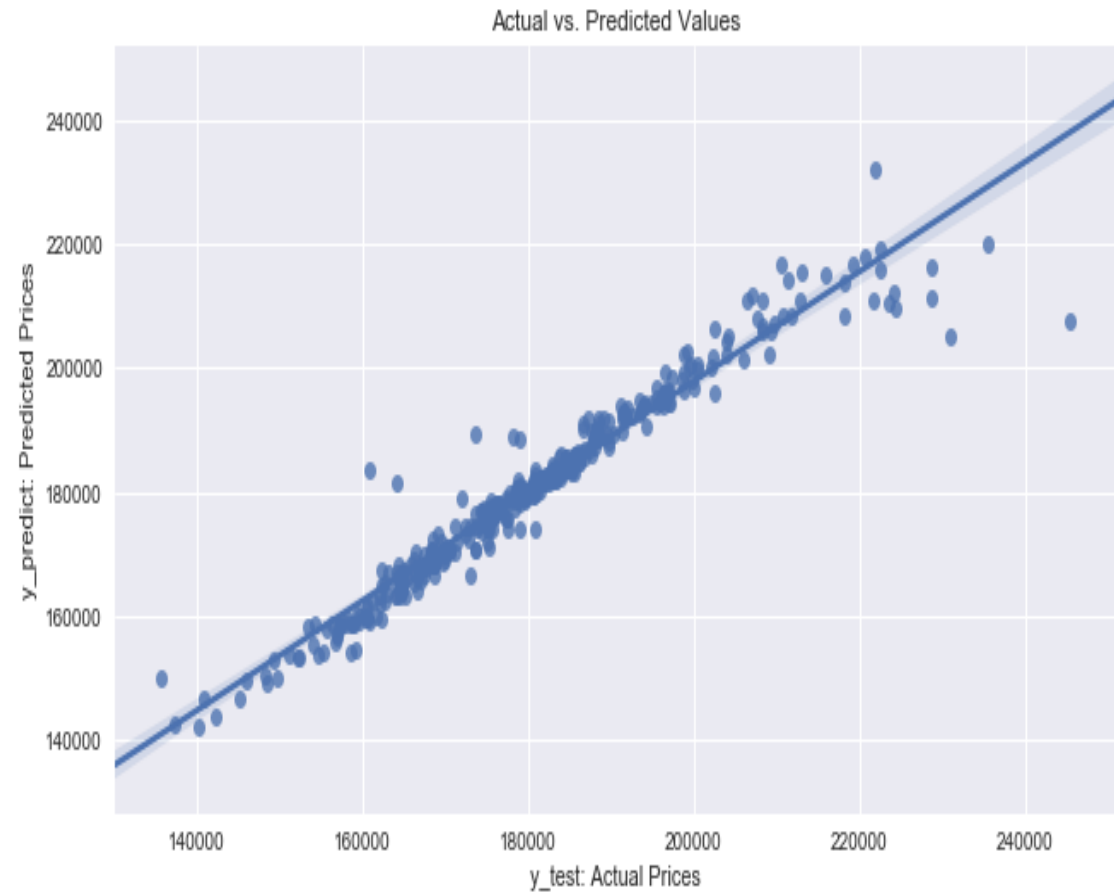
Results

Models Comparison:

Linear Regression RMSE: 42605.2337

Random Forest Regression RMSE: 30084.3291

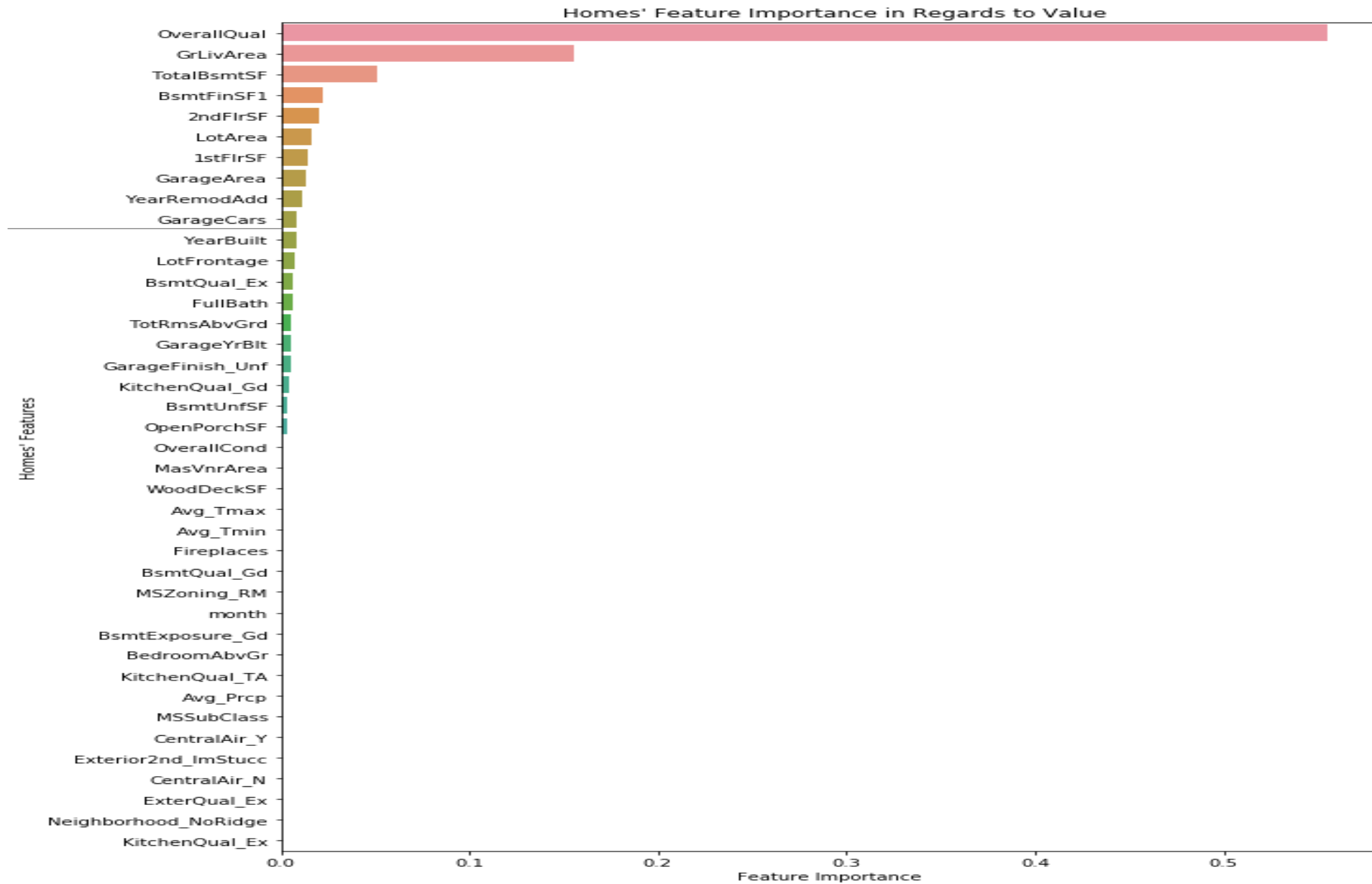
Gradient Boosting Regression RMSE: 27161.7518



Modeling and Analysis

FEATURE IMPORTANCE RELATED TO PRICE

	Features	Importance
3	OverallQual	0.555
15	GrLivArea	0.155
11	TotalBsmtSF	0.051
8	BsmtFinSF1	0.023
13	2ndFlrSF	0.020
2	LotArea	0.016
12	1stFlrSF	0.014
26	GarageArea	0.013
6	YearRemodAdd	0.011
25	GarageCars	0.009



Business Value

Sharper report analyses to help:

- Negotiation of sales
- Better Purchases
- Other strategic agreements related to real estate properties.