# Finding the Zestimate

## What is Data Science Department working on?

**Overview:**

⌂ What is the Zestimate and why do we need it?

⌂ How did we come up with the Zestimate?

⌂ How can we implement it?

⌂ Steps to further improve the Zestimate

# What is the Zestimate and why do we need it?

⌂ **Estimate**s the price of a house on **Z**illow.com

⌂ Important resource for customers

⌂ More listings, more revenue

⌂ More listings, more data

# How did we come up with the Zestimate?

⌂ Linear regression algorithm that can predict the prices of houses using historical data

⌂ **Ames Housing Data**:

- Residential properties sold in Ames, IA from 2006 to 2010
- 2980 rows (houses)
- 82 columns (features)

# How did we come up with the Zestimate?

⌂ **Data Cleaning**:

- ○ 26 feature columns with Null Values!!!
- ○ Dropped outliers

⌂ **Feature Engineering and Data Transformations**:

- ○ Log Transformation of sale price
- ○ 40 feature columns transformed to dummies or ordinal data
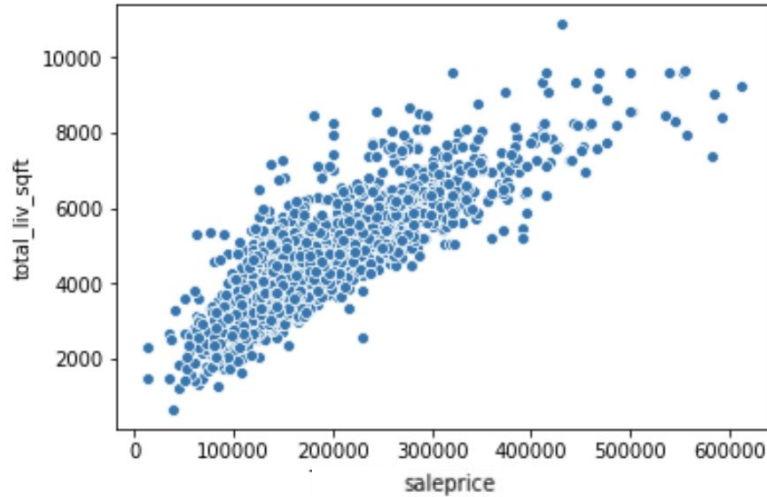- ○ Build new features out of correlated features

# Metrics

⌂ **$R^2$:** Percentage of variability in the data explained by Model

⌂ **Cross-Val-$R^2$ :** $R^2$ for 5 fold cross validation within the train data

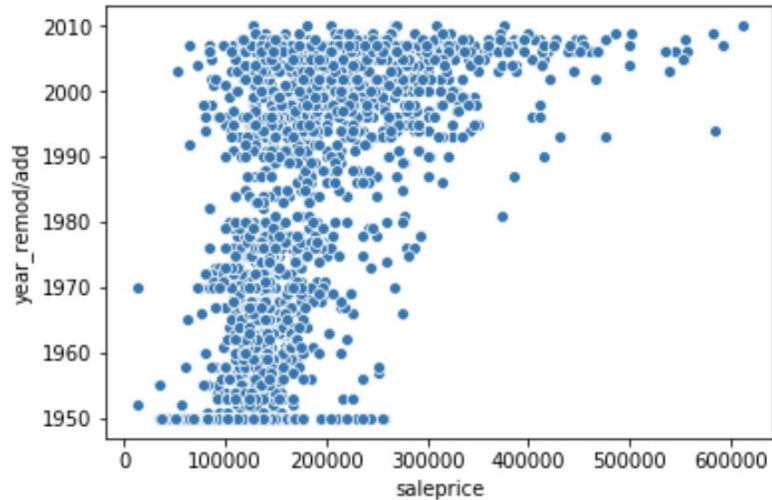⌂ **Mean Squared Error:** Mean of the squared residuals

# First Models: Linear Regression, no Regularization

⌂ **Feature Selection:**

Feature with a correlation coefficient with sale price



Total Living Sqft: corr 0.8



Year Remodeled: corr 0.55

# First Models: Linear Regression, no Regularization

| | Train R2 score | Cross val R2 score | Test R2 score | Mean Squared Error |
|---|---|---|---|---|
| OLS: corr >0.5 | 0.879156 | 0.870369 | 0.894068 | 0.017317 |
| OLS: all numeric | 0.922814 | 0.888454 | 0.907216 | 0.015168 |
| OLS: corr >0.4 | 0.881082 | 0.871384 | 0.893702 | 0.017377 |

**OLS corr > 0.5: 20 features**
**OLS corr > 0.4: 24 features**

**OLS all numeric: 106 features**
**→ Overfit!**

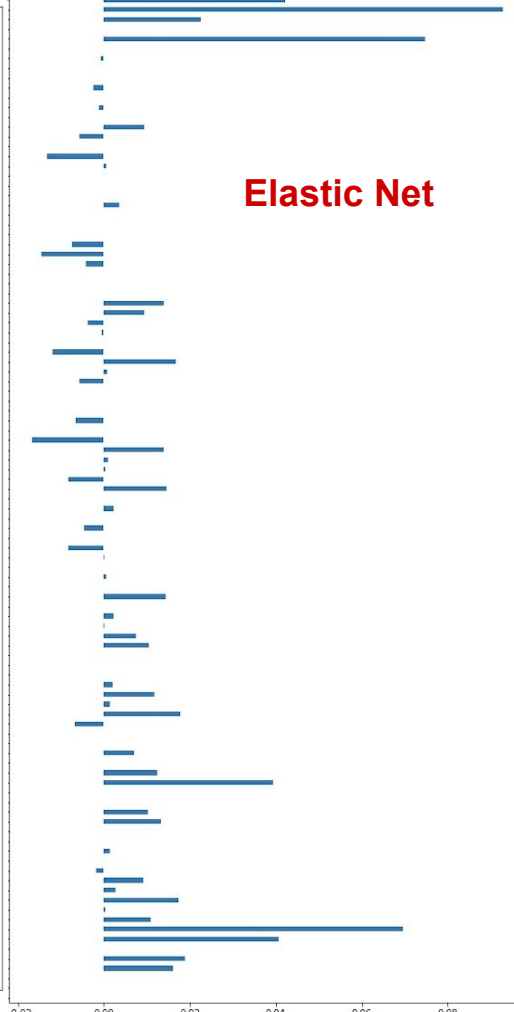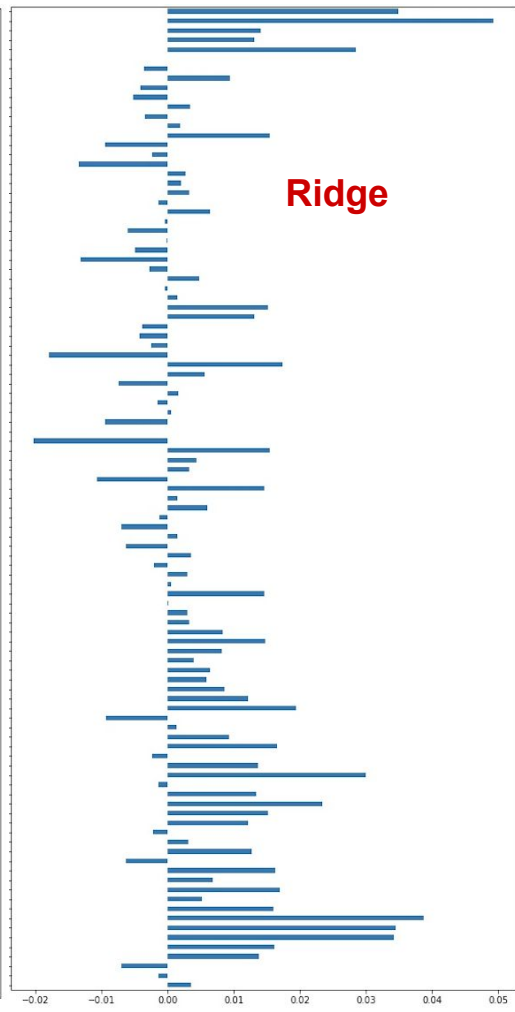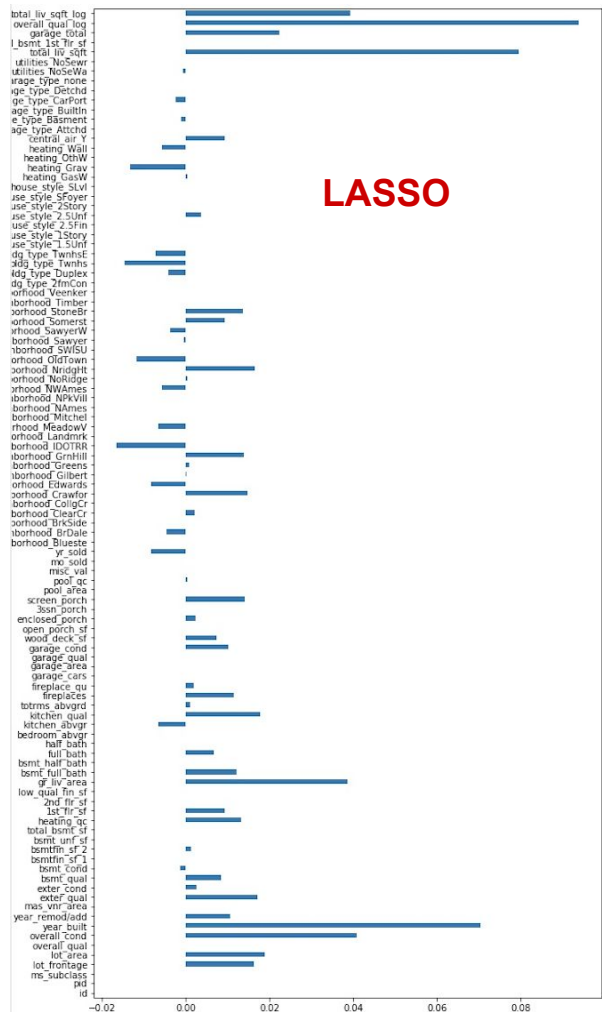# Improved Models: Linear Regression with Regularization

| | Train R2 score | Cross val R2 score | Test R2 score | Mean Squared Error |
|---|---|---|---|---|
| LASSO | 0.917852 | 0.901032 | 0.914756 | 0.013935 |
| Ridge | 0.918336 | 0.901632 | 0.913093 | 0.014207 |
| ElasticNet | 0.917990 | 0.901172 | 0.914818 | 0.013925 |

**Performance against external Test Data:**
**Ridge > Elastic Net > Lasso**

LASSO     Ridge     Elastic Net

# Next steps in the implement of the Zestimate

⌂ Right now: only applicable for Ames, IA

⌂ Future: Build similar model with the data from zillow.com

- Predictions for the whole US

- Include only features that are readily available

**Steps to further improve the Zestimate in the future:**

⌂ Experiment more with features and regularizations

⌂ Incorporate Location Data more strongly

⌂ Start a Data Science Competition on kaggle.com