



## Project 2 - Ames Housing Data

By Afolabi Cardoso



# Problem Statement

I am a Data Science Phd student at MIT and my Supervisor asked me to determine the best Linear Regression technique to use in determining the Sale Price of houses in his home town of Ames, Iowa.



# Overview



- Problem Statement
- Dataset
- OLS vs Ridge Regression vs Lasso
- Data Cleaning and Feature Engineering
- EDA
- Model Evaluation
- Recommendations and Conclusion



## Dataset

Location

**State Office  
Ames, Iowa**

Number of  
features

**79**

Samples

**2050**



## Data Cleaning and Feature Engineering

### Features Dropped

**'Lot Frontage'**  
**'Alley'**  
**Fireplace Qu'**  
**'Pool QC'**  
**'Fence','Misc Feature'**

### Top Features

**'Overall Qual'**  
**'Gr Liv Area'**  
**'Garage Area',**  
**'Total Bsmt SF'**  
**'1st Flr SF'**  
**'Year Built'**

### Dummied Features

**42**



# Models

OLS

Minimizes the MSE

LASSO

Penalized for the sum of absolute values of the weights

$$\frac{1}{2m} \sum_{i=1}^m (y - Xw)^2 + \alpha \sum_{j=1}^p |w_j|$$

RIDGE

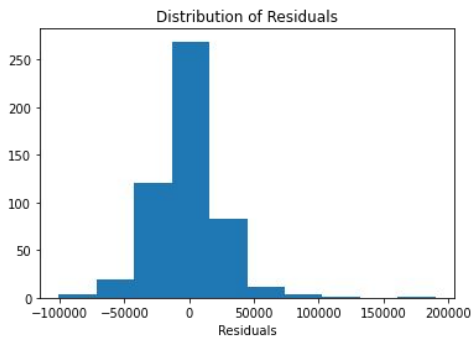
Penalizes the model for the sum of squared value of the weights

$$\sum_{i=1}^n (y - Xw)^2 + \alpha \sum_{j=1}^p w_j^2$$

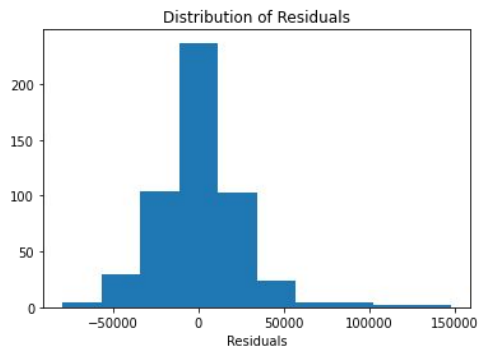
# Distribution



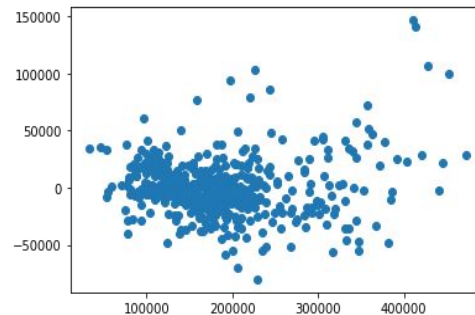
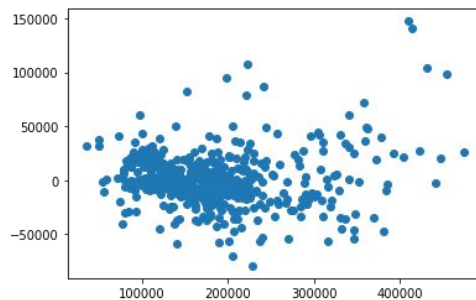
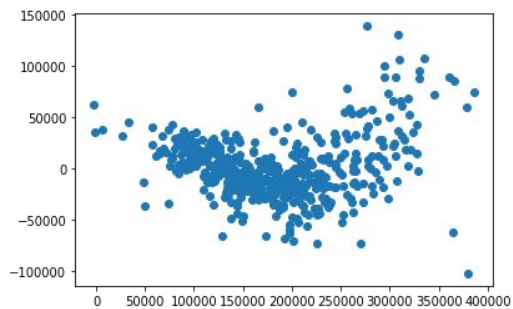
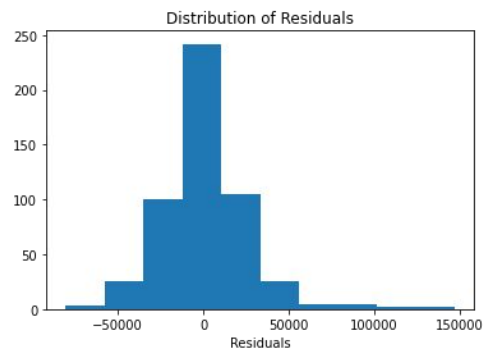
**OLS**



**Lasso**



**Ridge**



# Evaluation

	OLS (numerical)	OLS (dummied)	Lasso	Ridge
<b>Train set</b> R2	80%	89%	89%	89%
<b>Test set</b> R2	85%	88%	90%	91%
<b>Cross evaluation</b> R2	74%	79%	80%	82%
<b>RMSE</b>	30138	26437	25212	25024





## Conclusion

- The Ridge Model performed best overall
- The Lasso had similar results
- OLS performed better after dummifying the categorical features



## Recommendation

- I recommend to use either the Ridge or the Lasso model because they both gave a better overall result.
- As more data is collected the model will be reevaluated to see if the results still holds true.



“Data Science is an Art.”

—Riley

