# W207 Final Project

*AMES Housing*
*Mrinal Chawla, Thomas Gao, Fengyao Luo*

# The Inference Problem

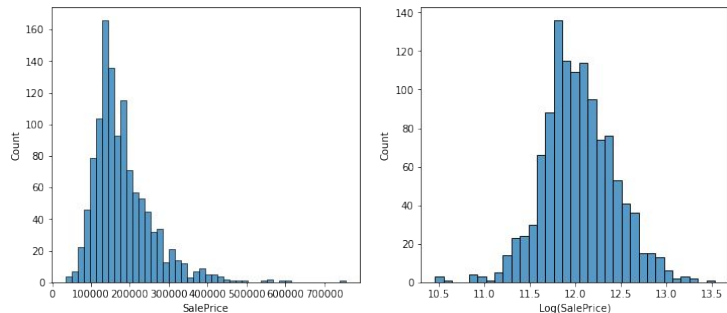| GIVEN |
| --- |
| <ul><li>A vector of features about the house<ul><li>Neighborhood</li><li>Quality<ul><li>Overall, Pool, etc.</li></ul></li><li>Sale Condition<ul><li>Normal, Abnormal</li></ul></li><li>Amenities<ul><li>Alley, Garage, Basement, etc.</li></ul></li><li>Size<ul><li>Sqft, # of rooms, # of bathrooms etc</li></ul></li></ul></li></ul> |

| PREDICT |
| --- |
| <ul><li>Sale Price of the House</li></ul> |

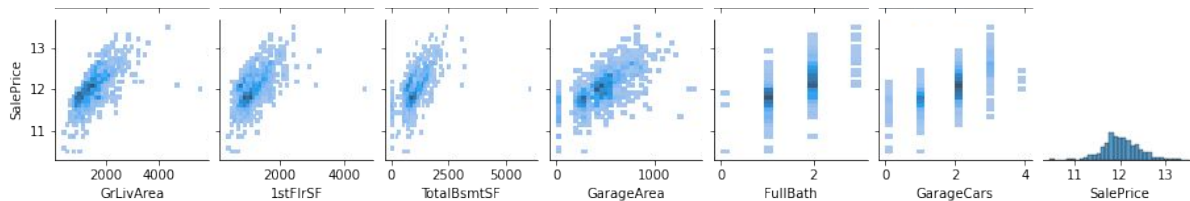| WHY |
| --- |
| <ul><li>Helps housing market for both sellers and buyers</li><li>Help plan renovations</li><li>Help plan infrastructure improvements</li><li>Etc.</li></ul> |

# Exploratory Data Analysis

- SalePrice slightly right-skewed
  - Log transformation to fix
- Sizing variables highly correlated
  - Bsmt sqft vs 1st floor sqft
  - GarageArea vs GarageCars

# Exploratory Data Analysis

- Generally linear relationships

- Few outliers for expensive homes

- Top indicators

  - Overall Quality

  - Living Area

  - Neighborhood

- Scaling and encoding

# Baseline Models

PREDICT MEAN

Predict average sale price for every house

RMSE: 0.419

LINEAR REGRESSION

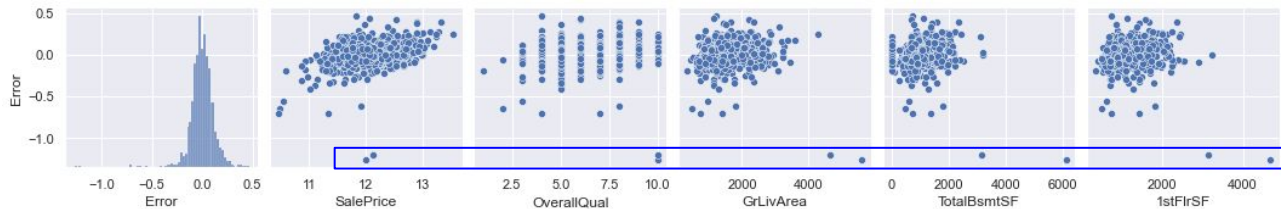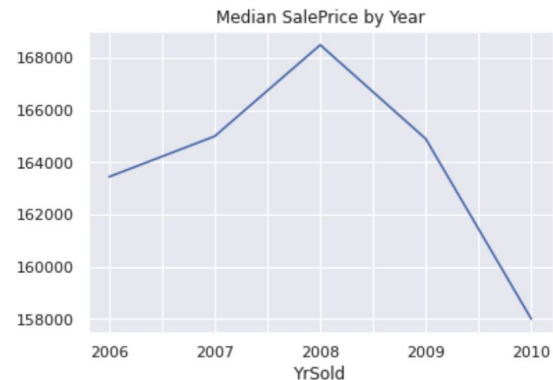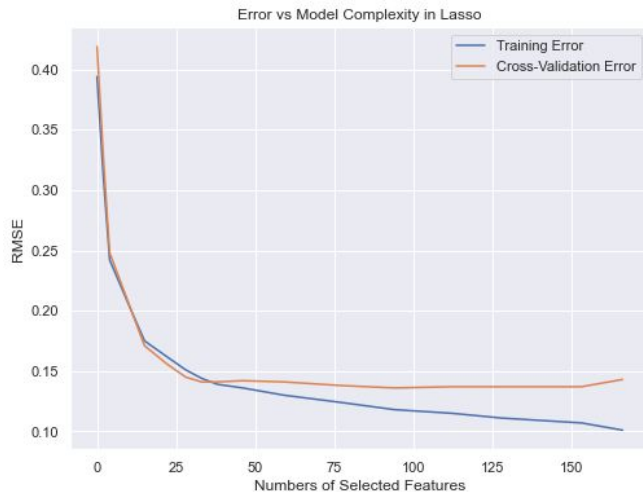Use top 2 features and neighborhood

RMSE: 0.169

LASSO

Linear Regression with L1 Regularization
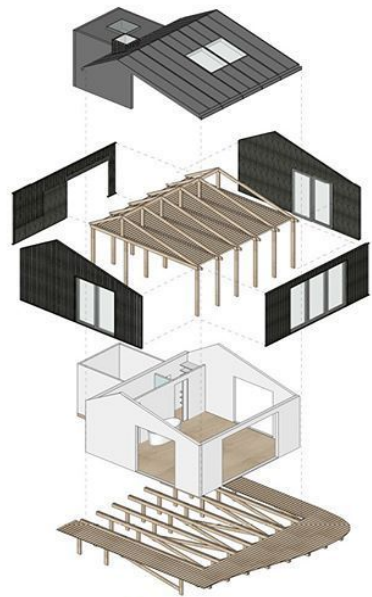
RMSE: 0.137

# Error Analysis

- Complexity vs Performance
  - Number of Features
  - RMSE
- Outliers
  - Living Area
  - Quality
- Non-Linear Relationship
  - YrSold
  - SalePrice

# Feature Engineering

- Aggregate size features
  - Total Sqft
  - Average Room Sqft
  - Total bathrooms
  - Total porch sqft
- Presence of amenities (binary)
  - Alley
  - Garage
  - Basement
  - Pool
- Years Since Remodelled  (Year Sold - Year Remodelled)
- Seasonality (Month Sold → Season)
- Skewness (np.log)
- Neighborhood bins

# Final Models

| LASSO | SPLIT LASSO |
|---|---|
| Linear Regression with L1 Regularization | Linear Regression with L1 Regularization + Separate Model for each Neighborhood Bin |
| RMSE: 0.121 | RMSE: 0.126 |
| RANDOM FOREST | ENSEMBLE |
| Nonlinear Model | 0.75 Lasso + 0.25 Random Forest |
| RMSE: 0.135 | RMSE: 0.119 |

# Final Performance

- Model: Ensemble

- RMSE: 0.1217

- Kaggle: 10th percentile