
Housing Price Analysis

Team 2: Mckinlytic
Karl Hu & Ellen Wang

Table of Contents

01 Background

02 Data
Preparation

03 Descriptive
Statistics

04 Regression
Analysis

05 Classification
Analysis

06 DataRobot
Comparison

07 Conclusions &
Insights

Background

- ❖ Volume of Real Estate Market
- ❖ Preliminary Screening
- ❖ Identify Key Taste Preference
- ❖ Time saving and efficiency



Data Preparation

- ❖ Quality
- ❖ Transformation
- ❖ Normalization
- ❖ Partition/Validation



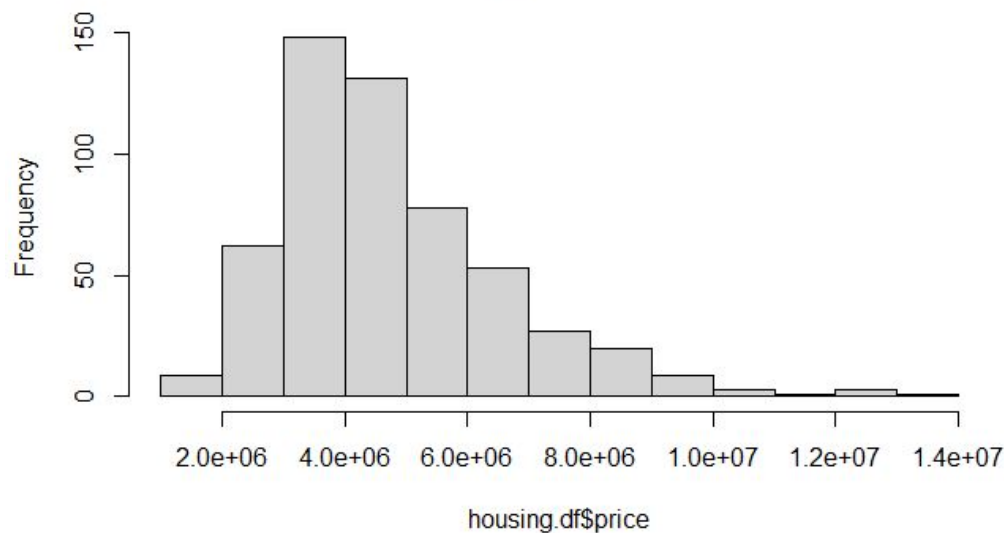
Independent Variables

- | | |
|----------------------------------|--------------------------------------|
| 1. Total area of the property | 2. Number of bedrooms |
| 3. Number of bathrooms | 4. Number of stories/floors |
| 5. Number of parking spaces | 6. Access to the main road |
| 7. Availability of a guest room | 8. Availability of a basement |
| 9. Availability of water heaters | 10. Availability of air conditioning |
| 11. Postal preferred area status | 12. Furnishing status |

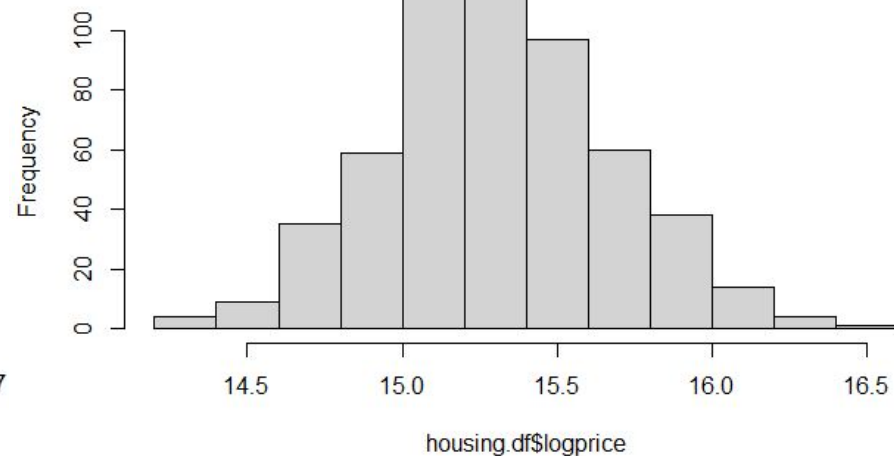
Dependent Variable: Log_Price (Regression), Price Level (Classification)

Descriptive Statistics

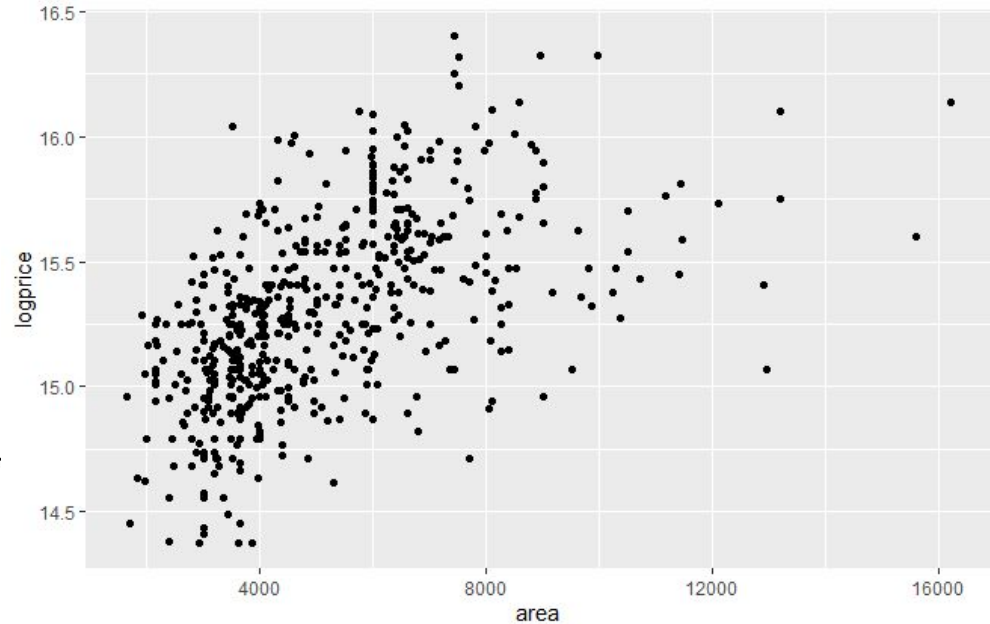
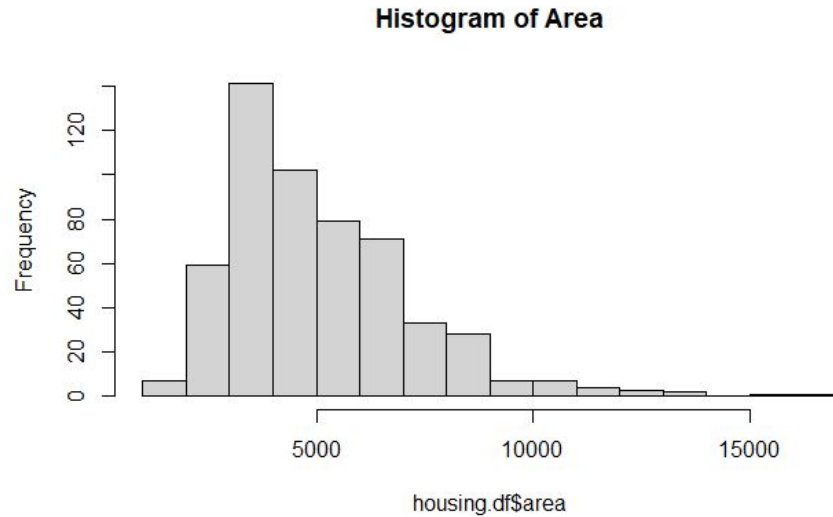
Histogram of Price



Histogram of Log of Price



Descriptive Statistics - cont.



Regression Analysis

Multiple Regression

0.6907

Best Model

Regression Tree

0.5946

Penalty: 0
Min-Leaf: 21

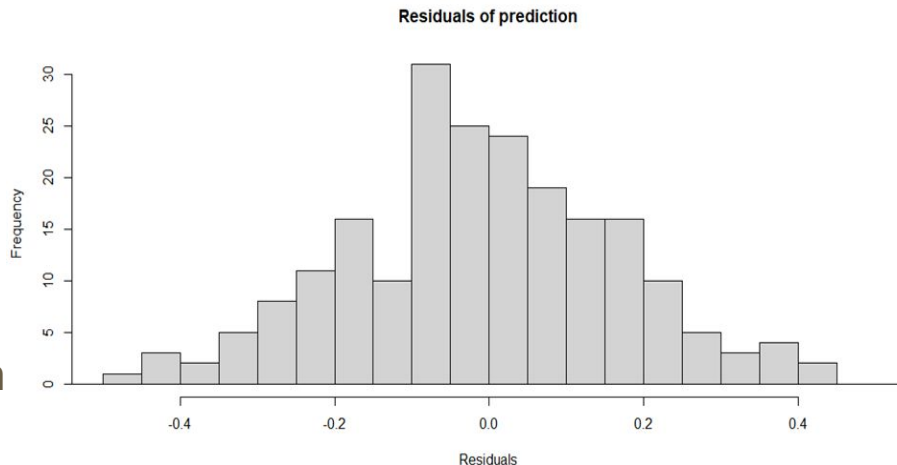
Random Forest in R

0.6698

Memory Shortage

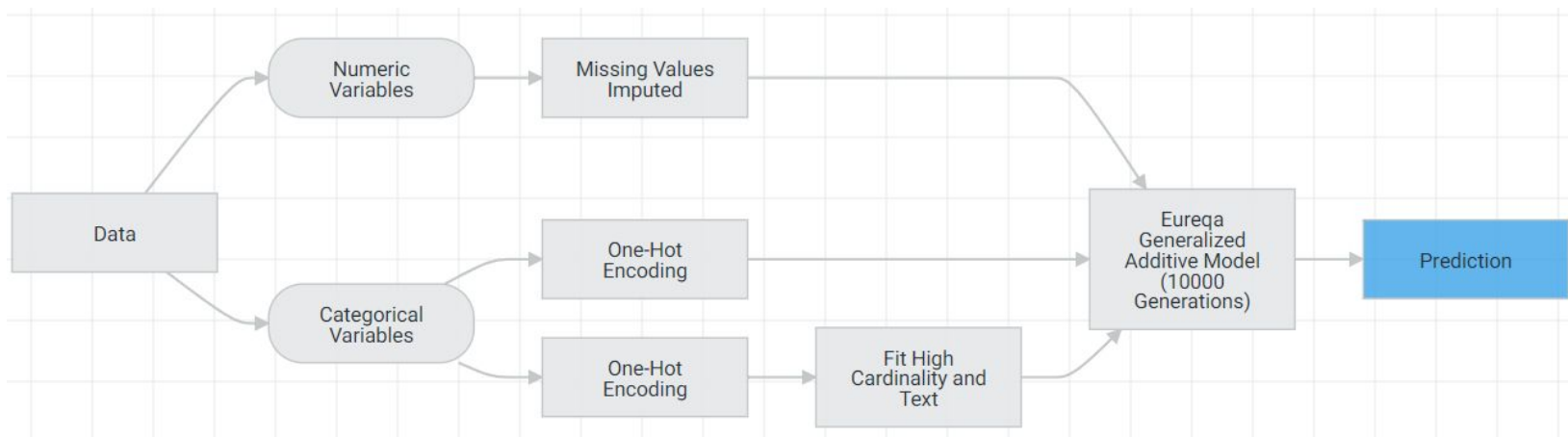
Multiple Regression - cont.

- ❖ Started with simple multilinear regression
- ❖ Removed insignificant variables
- ❖ Removed 7 outliers
- ❖ Feature Engineering and Interaction Terms
- ❖ Cluster Variables
- ❖ RMSE 0.2069
- ❖ R-Squared 0.6907



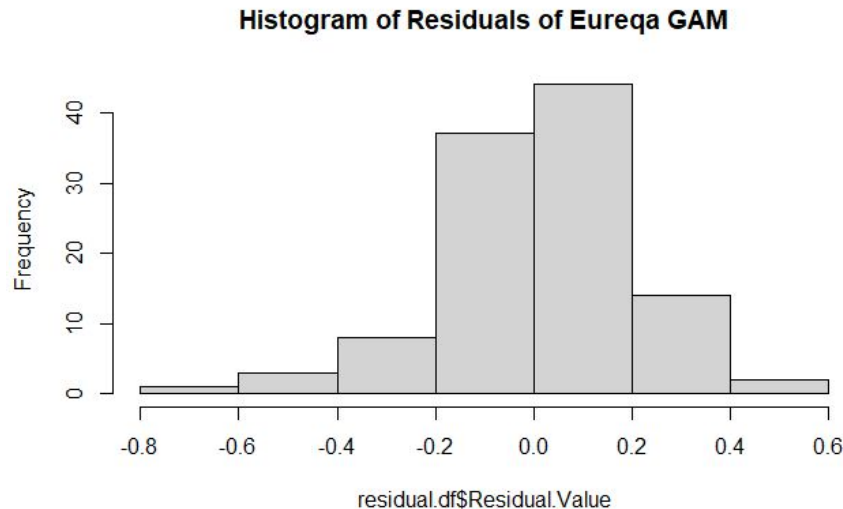
Datarobot - Best Model

Eureqa Generalized Additive Models (Eureqa GAM) - 10000 Generations



Eureqa GAM - cont

- ❖ Unique partition logic
- ❖ Complexity: 14
- ❖ Error: 0.010(Surrogate MSE)
- ❖ RMSE 0.1986
- ❖ R-Squared 0.6627



Regression - DataRobot Comparison

DR vs Human

Multilinear Regression

Eureqa GAM

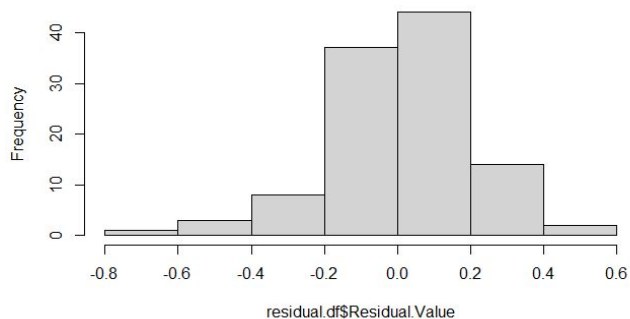
Feature List

Removed 5 variables; Feature Engineering
Cluster Variables

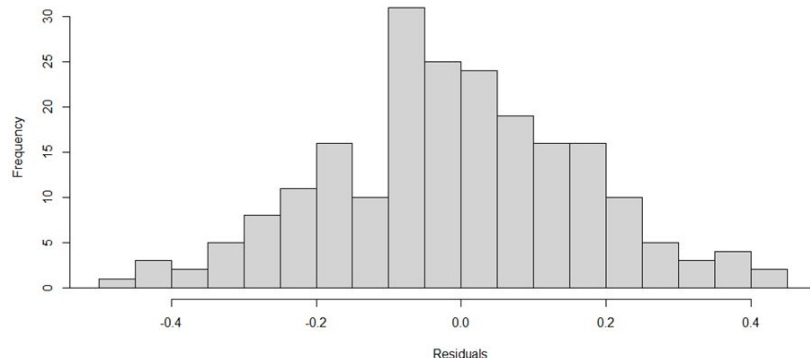
All Features

Residuals

Histogram of Residuals of Eureqa GAM



Residuals of Multilinear Regression



R-Squared

0.6907

0.6627

RMSE

0.2069

0.1986

Classification - Non-Ensemble Models

Logistic Regression*

0.9161

Target Variable

- Price level

KNN

0.8546

Grid Search

- KNN: number of k
- Tree: cp, minsplit

Classification Tree

0.8410

Best Model

- Best non-ensemble: logistic regression

Best Hand-Crafted Model - Logistic Regression

- ❖ Started with simple logistic regression on all the variables (excluding price)
 - Removed price to avoid multicollinearity
- ❖ Regress on cluster variables
- ❖ Split data into training set (60%) and validation set (40%)
- ❖ Removed insignificant variables
 - Furnishing status, basement, parking, mainroad, bedrooms
- ❖ Removed 19 outliers
 - Greatly improved the accuracy of the model
 - Also suggested by DataRobot
- ❖ Checked confusion matrix
 - Overall accuracy = .9221, F1 score = .9161

Classification - Ensemble Models

Random Forest

0.8629

Hyperparameters

- Boos
- Mfinal
- Coeflearn (for Boosted)

Bagged

0.8183

Value Range for
Grid Search

- [TRUE, FALSE]
- [1,10]
- ['Breiman', 'Zhu', 'Freund']

Boosted

0.8653

Optimal Values

- Boos = FALSE
- Mfinal = 8
- Coeflearn = 'Breiman'

Classification - DataRobot Comparison

	Best Hand-Crafted Model	Best DataRobot Model
Model of Choice	Logistic Regression*	Light Gradient Boosting Tree
Feature List	Removed 6 variables; Cluster variables	All Features (excluding price)
Top Factors By Importance	<div>1. Air conditioning</div> <div>2. Hot water heating</div> <div>3. Guestroom</div>	<div>1. Area</div> <div>2. Stories</div> <div>3. Bathrooms</div>
Confusion Matrix	<div><div>172</div><div>14</div><div>10</div><div>112</div></div>	<div><div>47</div><div>17</div><div>1</div><div>44</div></div>
AUC	0.9831	0.9372

Conclusions & Insights

❖ Model Building

- Less is more
 - Insignificant variables
 - Outliers
- Tradeoffs between interpretability vs accuracy

❖ Results

- Which variables are most important in predicting house price:
 - area
 - Bathrooms
 - airconditioning

Thank you!