

Residential Assessment Neighborhoods

05-Feb-09

(#) Name [Multiplier]

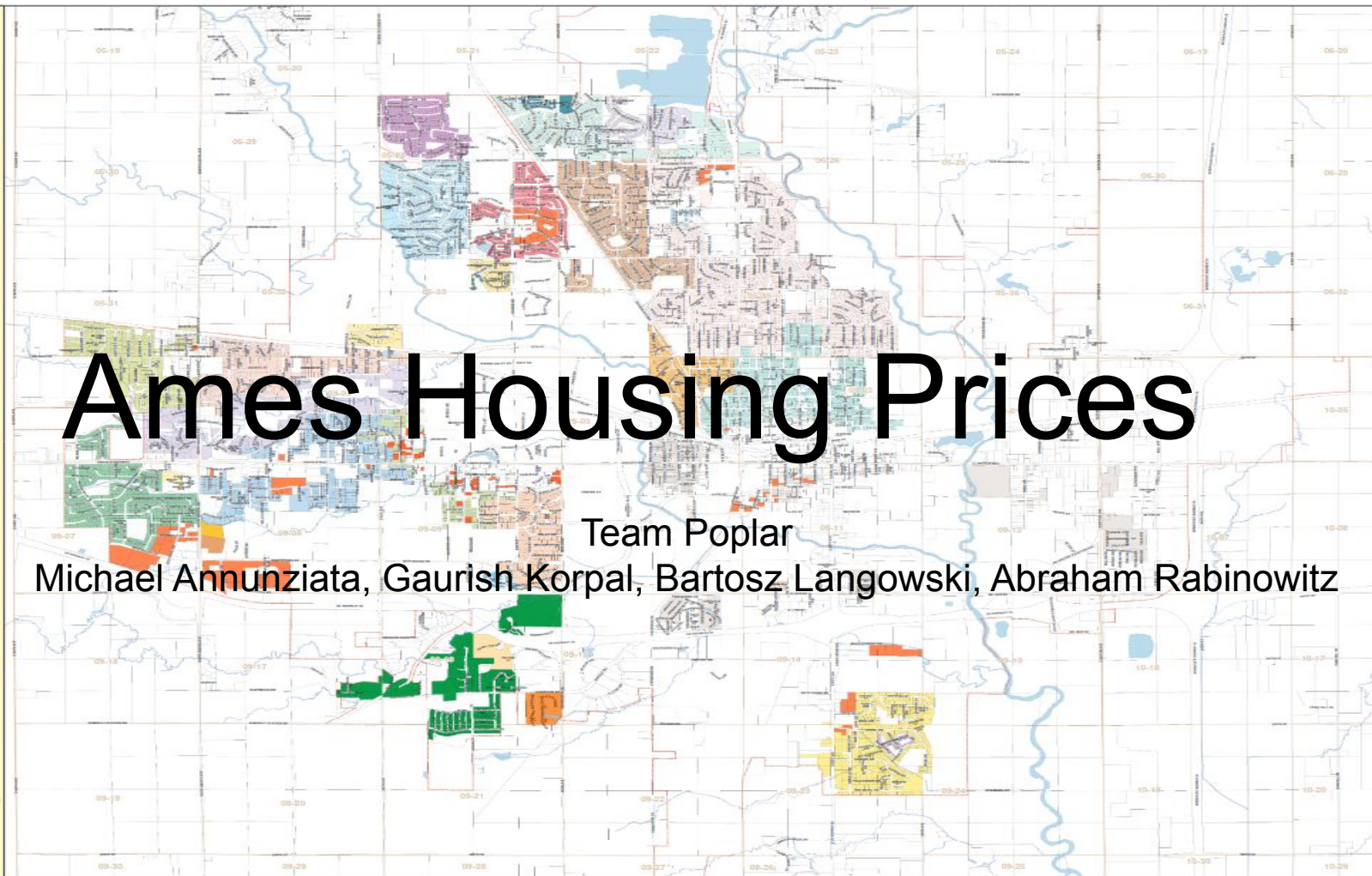
- 1 (8) IOC ISU [54]
- 1 (9) IOCCons [89]
- 1 (10) MacCondo [99]
- 1 (11) Br/Dale [102]
- 1 (12) NoPKW [109]
- 1 (13) Greens [113]
- 1 (14) MeadowV [30]
- 1 (15) Bluestm [99]
- 1 (16) WilwCr1 [81]
- 1 (17) WilwCr2 [85]
- 1 (18) LandmK [54]
- 1 (19) GrmHil [154]
- 1 (20) Stonebr [104]
- 1 (21) Blmngtn [105]
- 1 (22) NPrdgH [104]
- 1 (23) Wessex [93]
- 1 (24) NoRidge [101]
- 1 (25) Veswke [98]
- 1 (26) Timber [103]
- 1 (27) ClearCr [103]
- 1 (28) Somerst [101]
- 1 (29) Gilbert [97]
- 1 (30) NW Ames [99]
- 1 (31) N Ames [100]
- 1 (32) BrkSide [106]
- 1 (33) OldTown [102]
- 1 (34) IDOT&RR [102]
- 1 (35) Mitchel [99]
- 1 (36) Crawfor [106]
- 1 (37) S&W ISU [99]
- 1 (38) Edwards [98]
- 1 (39) Sawyer [101]
- 1 (40) Sawyer [98]
- 1 (41) CollgCr [98]

City of Ames
Assessor's
Office

E



1:13,200



Ames Housing Prices

Team Poplar

Michael Annunziata, Gaurish Korpai, Bartosz Langowski, Abraham Rabinowitz

The Problem

1. Accurately predicting housing prices in Ames, Iowa.
2. Understanding what features of a home influence its value.

The Dataset

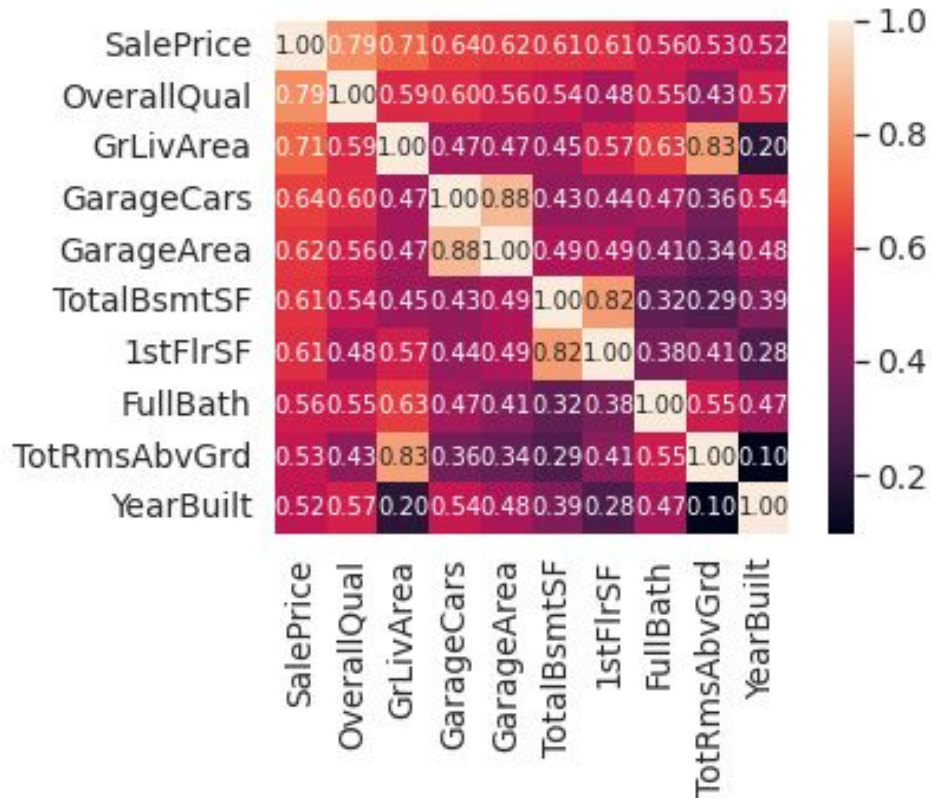
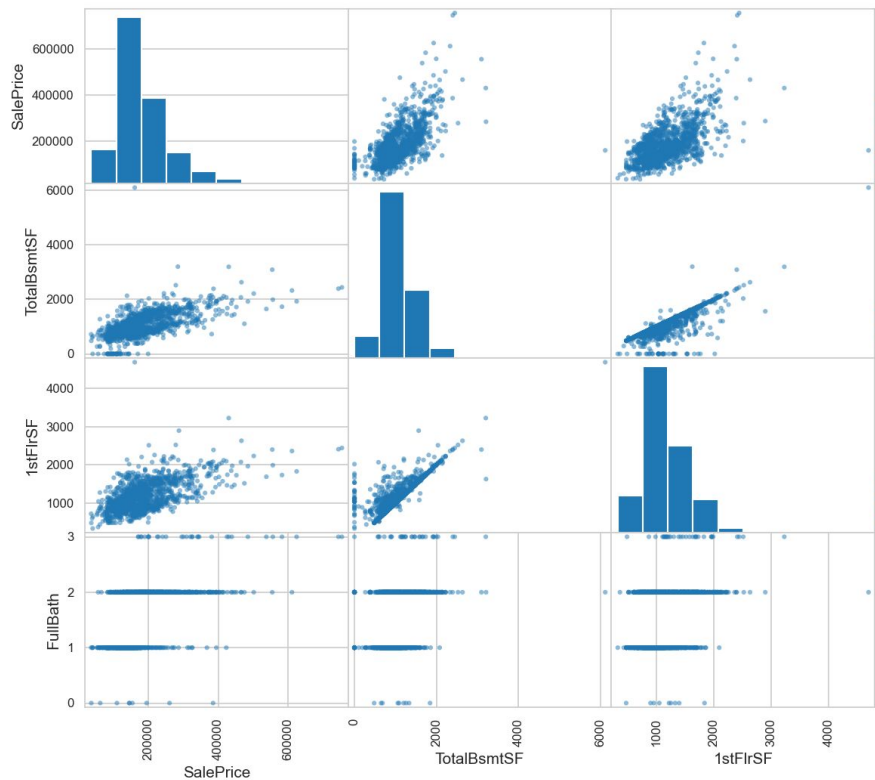
Dean De Cock (2011) Ames, Iowa: Alternative to the Boston Housing Data as an End of Semester Regression Project, Journal of Statistics Education, 19:3, DOI: [10.1080/10691898.2011.11889627](https://doi.org/10.1080/10691898.2011.11889627)

Anna Montoya (2016), House Prices - Advanced Regression Techniques. Kaggle. <https://kaggle.com/competitions/house-prices-advanced-regression-techniques>

- Provides a description of the sale of 2930 individual residential property in Ames, Iowa from 2006 to 2010.
- Contains 80 explanatory variables with numerical and categorical data such as square footage, neighborhood, and overall house condition.
- The training data consists of 1460 observations.

Exploratory Data Analysis

- To study correlation between the **Sale Price** of the house and numerical features we used:
 - *Correlation matrix*
 - *Heat map*
 - *Scatter plots*



Modeling Approach

- Linear Regression:
 - Took advantage of a strong linear relationship between sale price and numerical variables.
 - Elastic Net combined Lasso and Ridge regressions:
 - Performed better than ordinary least squares, lasso, or ridge regression alone.
 - Did not incorporate information from many categorical variables.

Modeling Approach

- Ensemble Learning:
 - Considered several models to better capture nonlinearity and incorporate categorical variables.
 - Final ensemble model is a stacked regression combining Elastic Net, XGBoost, and LightGBM.
 - Cross validation to tune hyperparameters. For example:
 - $\alpha = 0.0006$
 - L_1 ratio = 0.9
 - Learning rate = 0.05
 - n estimators = 1000

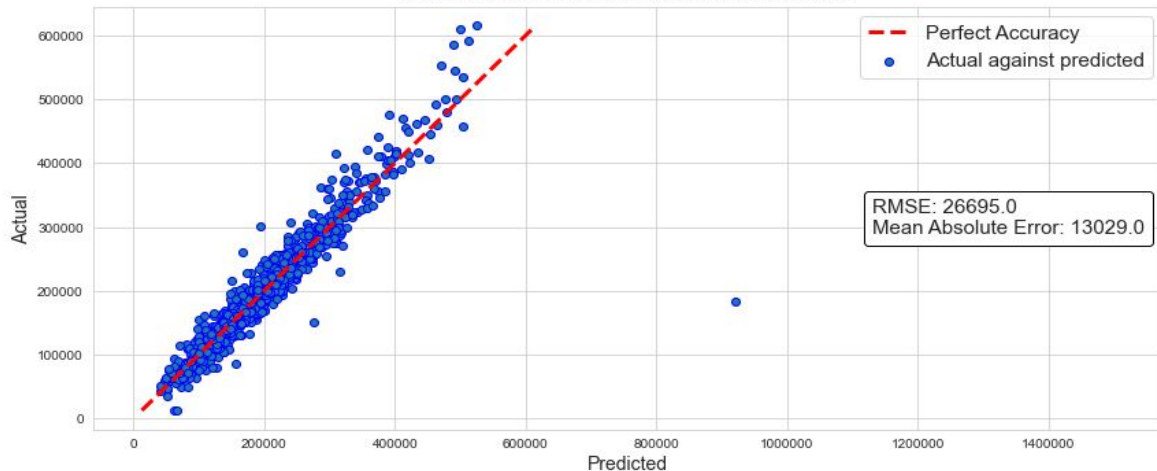
Modeling Approach

- Used cross validation to select final model:

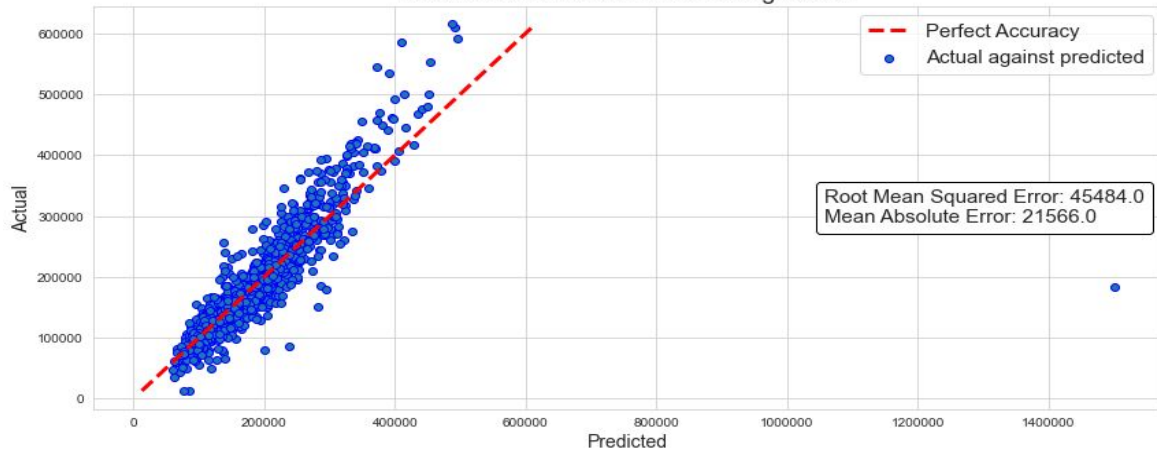
Model	Log Mean Squared Error
Ridge Regression	.0125
Elastic Net Regression	.0122
LightGBM	.0142
XGBoost	.0136
Ensemble (Elastic Net, LGBM, XGBoost)	.0117

Results

Test Performance of the ensembled model



Test Performance of Baseline Regression



Performance on Unseen Data:

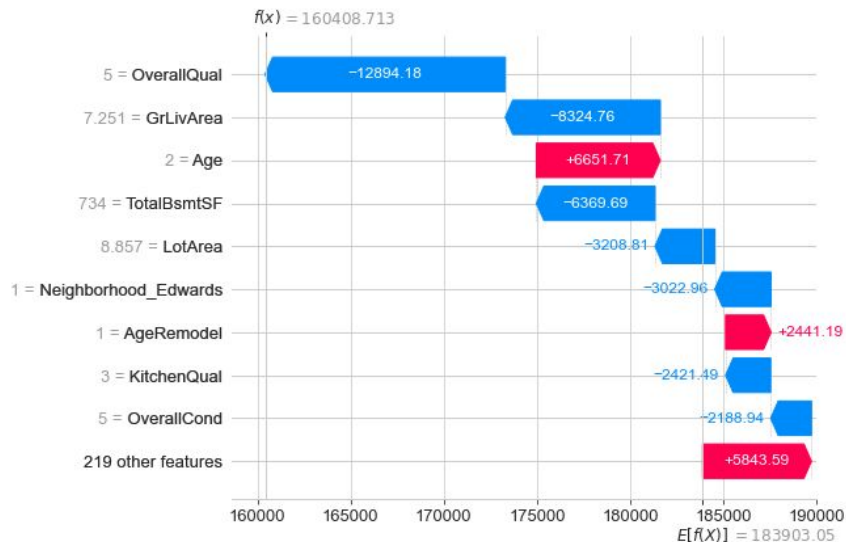
- More accurate than a baseline linear model.
- Better at pricing more expensive houses.
- Cannot handle extreme outliers.
- Further work can focus on handling such outliers or enlisting domain experts to help price them.

Results

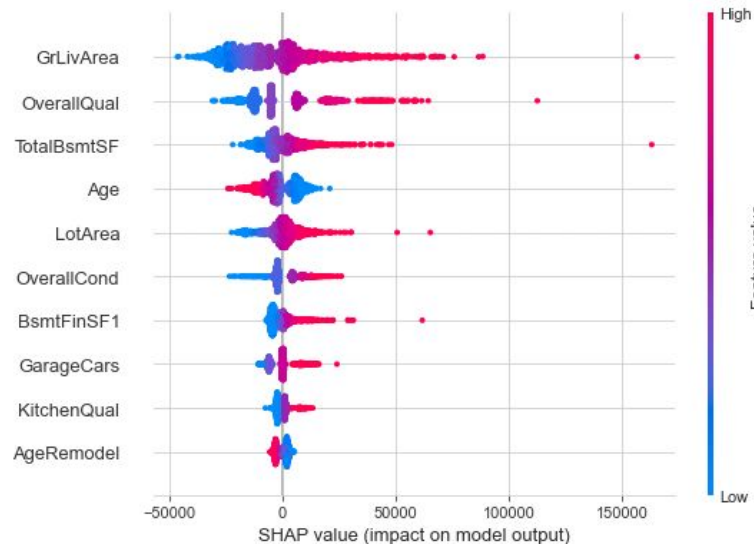
Interpreting our model with SHAP:

- SHAP values quantify the effect that changing values of specific features has on the way that the model predicts Sale Price.

SHAP values for a sample prediction



Distribution Summary of SHAP Scores



Appendix: Exploratory Data Analysis

Analysis of categorical features:

- Performed *one-hot encoding*.
- We used *violin plots* to understand the feature importance:

