

STAT 542 Project 1: House Price Prediction

Wenbo Fu (679744457), Bingyan Liu (668046518)

1 Data and Objective

The data contains the training data and the test data, with the sale price as the target prediction feature. The goal is to reduce the Root-Mean-Squared-Error (RMSE) between the logarithm of the predicted price and the logarithm of the observed sales price.

$$\sqrt{\frac{1}{n.test} \sum_{j=1}^{n.test} (\hat{y}_j - y_j)^2} \quad (1)$$

The RMSEs should be less than 0.125 for the initial 5 folders and less than 0.135 for the last 5 folders.

The details can be found in https://liangfgithub.github.io/Proj/F23_Proj1.pdf.

2 Method

Our algorithm contains two parts: data preprocessing and model fit.

The data preprocessing part takes the training data and testing data as inputs, which contains the following steps:

- On the training data set, fill the missing values in the Garage_Yr_Blt column with value 0.
- Remove imbalanced features, where imbalanced features are determined by more than 95% of the sample feature values are the same.
- Winsorize numerical features, by setting all values higher than 0.95 quantile of all sample feature values to the 0.95 quantile.
- Standardize the numerical features with the StandardScaler class.
- One hot encode categorical features with the OneHotEncoder class.
- On the testing data set, follow similar steps: fill the Garage_Yr_Blt column with value 0; remove imbalanced features found from the training data; winsorize numerical

features with the 0.95 quantile values computed from the training data; standarize the numerical features with the fitted StandardScaler from the training data; one hot encode the categorical features with the fitted OneHotEncoder from the training data.

- Output the processed traning data and testing data.

The model fit part takes the processed traning data and the testing data as inputs, which contains the following steps:

- In the linear model part, we choose elastic net model with l_1 ratio 0.5, the best λ is searched with the grid search and 5 fold cross validation, the candidate values are from the set $\{10^{-4+0.5(i-1)}\}$ with i integers from 1 to 15.
- In the tree model part, we choose XGBoost regressor with 5000 trees, 0.05 learning rate, max depth 6, subsampling rate 0.5.
- Use the processed training data to fit the two models and use the fitted model to fit the processing test data to give two predictions.
- Save the two predictions in mysubmission1.txt and mysubmission2.txt.

3 Result and Discussion

The linear model RMSEs on the 10 folds:

$$0.1204, 0.1167, 0.1169, 0.1119, 0.1100, 0.1344, 0.1346, 0.1317, 0.1342, 0.1257 \quad (2)$$

The XGBoost model RMSEs on the 10 folds:

$$0.1173, 0.1202, 0.1121, 0.1187, 0.1082, 0.1263, 0.1337, 0.1258, 0.1342, 0.1172 \quad (3)$$

The RMSEs satisfy the criterions. The XGBoost model result is slightly better (average RMSE 0.1214) than the linear model restul (average RMSE 0.1244). The whole computation time is about 115 seconds on a Macbook Air, 8GB memory, Apple M1 chip.

We did not work too much on hyperparameter tuning. We initially used methods described in the instructor's posts: <https://campuswire.com/c/G06C55090/feed/212>, <https://campuswire.com/c/G06C55090/feed/213>, but the RMSE from fold 6 is slightly greater than 0.135. Thus we improved the preprocessing process without pre-specifying feature columns to transform as in the previous posts, which give similar RMSEs but the criterions are satified for all 10 folds.