

HOUSE PRICES-ADVANCED REGRESSION TECHNIQUES

Tao Wang

University of Chinese Academy of Sciences, China

Introduction

For home buyers, they generally do not buy homes with basements or near rail-roads, and there are other features of a home that can even much more influence the price of a home than the number of bedrooms. This project provides 79 characteristics of a house that are used to predict the price of a house.

- Project Evaluation There are 1459 data in the test set, and the output contains ID numbers and predicted house prices. Submissions are evaluated on Root-Mean-Squared-Error (RMSE) between the logarithm of the predicted value and the logarithm of the observed sales price.
- File Description
 - train.csv – the training set
 - test.csv – the test set
 - data`description.txt – full description of each column, originally prepared by Dean De Cock but lightly edited to match the column names used here
 - sample`submission.csv – a benchmark submission from a linear regression on year and month of sale, lot square footage, and number of bedrooms

Data Analysis

There are 79 variables of inconsistent types, some discrete and some continuous, and after data processing, the variables of type object can be seen as follows:

#	Column	Non-Null Count	Dtype
0	<i>MSZoning</i>	1460	<i>non – null object</i>
1	<i>Street</i>	1460	<i>non – null object</i>
2	<i>Alley</i>	91	<i>non – null object</i>
3	<i>LotShape</i>	1460	<i>non – null object</i>
...	<i>object</i>

The characteristics of variables of type object are represented by string,such as:

Mszoning:Identifies the general zoning classification of the sale	
<i>A</i>	<i>Agriculture</i>
<i>C</i>	<i>Commercial</i>
...	...

The numeric variables are as follows:

#	Column	Non-Null Count	Dtype
0	<i>MSZoning</i>	1460	<i>non – null int64</i>
1	<i>Street</i>	1460	<i>non – null int64</i>
2	<i>Alley</i>	1201	<i>non – null float64</i>
3	<i>LotShape</i>	1460	<i>non – null int64</i>
...

Conclusion

The difference between the house price prediction experiment and the usual data classification is that the data classification only needs to give the data a specific prediction result and then analyze the correct rate. This time, however, the house price prediction is to give the prediction result to a more detailed point, and then to analyze the error. I think this kind of prediction is more difficult to go up to a better result, and the number of data sets itself is not large. I think the training effect can be improved by increasing the dataset through data expansion. In conclusion, I learned a lot about the new training model in this experiment. For the xgboost model, more adjustments are needed to increase the number of parameters and the depth of the tree to prevent underfitting.

Data Pre-processing

First we count the missing data in the training and test sets, and the statistics are as follows:

- Number of missing data in the training set:

<i>type</i>	num	persent
<i>PoolQC</i>	1453	0.995
<i>MiscFecture</i>	1406	0.963
<i>Alley</i>	1369	0.938
<i>Fence</i>	1179	0.808
<i>FireplaceQu</i>	690	0.472
<i>LotFronttage</i>	259	0.177
<i>GarageType</i>	81	0.055
<i>GarageYrBlt</i>	81	0.555
<i>GarageFinish</i>	81	0.555
...

- Number of missing data in the test set:

<i>type</i>	num	persent
<i>PoolQC</i>	1456	0.997
<i>MiscFecture</i>	1408	0.964
<i>Alley</i>	1352	0.926
<i>Fence</i>	1169	0.800
<i>FireplaceQu</i>	730	0.5
<i>LotFronttage</i>	227	0.155
<i>GarageYrBlt</i>	78	0.053
<i>GarageFinish</i>	78	0.053
...

- Remove the feature PoolQC, MiscFeature, Alley,Fence and FireplaceQu
- Use the average to fill in 8 missing numbers of masvnrarea, and use the mode to fill in the missing data of LotFrontage and GarageYrBlt.
- Use log transformation to reduce skewness.
- Delete unreasonable data
- Use mode to fill in the missing data of string type.

Training Result

- Training Result Summary

<i>method</i>	result
<i>LinearRegression(low_corrdeleted)</i>	0.7313
<i>RandomForest(low_corrdeleted)</i>	0.16099
<i>K – neighbor(low_corrdeleted)</i>	0.22111
<i>AdaboostRegressor(low_corrdeleted)</i>	0.22124
<i>Xgboost(low_corrdeleted)</i>	0.0.60098
<i>Xgboost</i>	0.17429
<i>GradientBoostingRegressor(low_corrdeleted)</i>	0.15922
<i>GradientBoostingRegressor</i>	0.14173