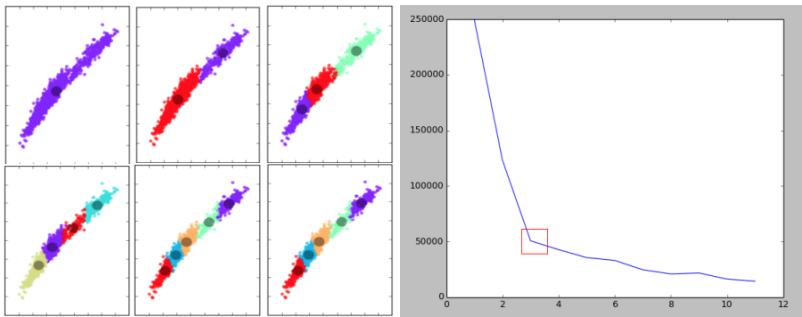


Introduction

In the real world, the house price is really a hot topic around the world. There are lots of features which can have impact on the sale price of a house. Hence, our team consider it is a good field to do data mining and analysis in order to guide real estate agents during decision making and pricing. Also, fundamentally, we can help house buyers to find the best house according to our analysis. Even some big companies like Airbnb can use the prediction of the house prices to inform their hosts to better pricing their places.

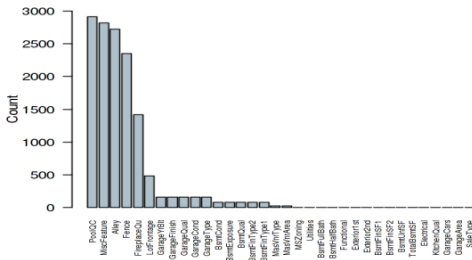
After discussion, our team choose the "House Prices: Advanced Regression Techniques" as our dataset, which is a good representative dataset from Kaggle platform.

Data Preprocess



Missing Data

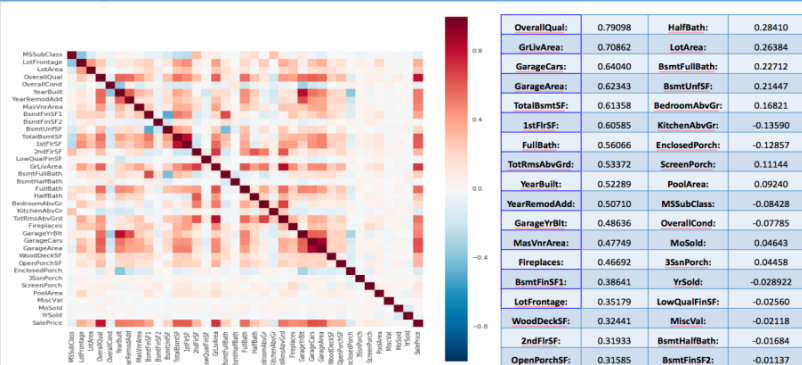
34 variables with missing values in dataset



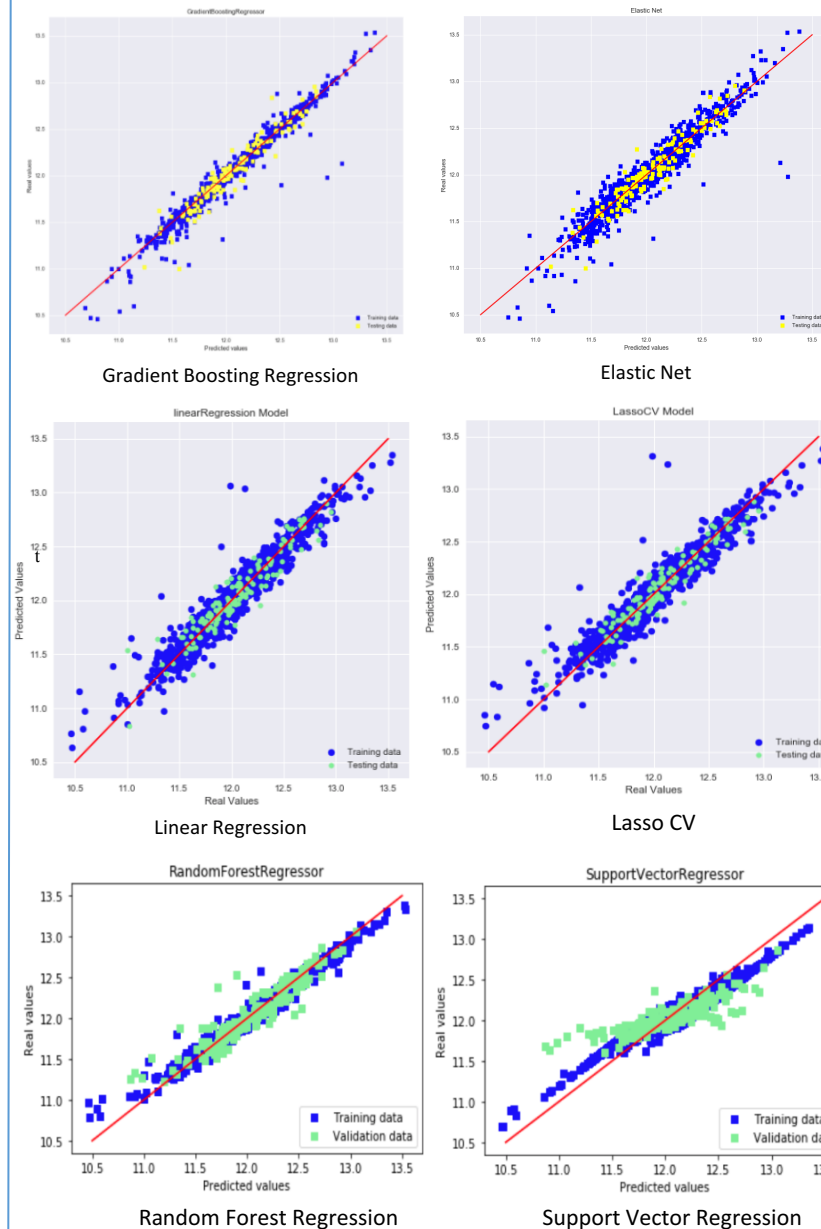
The missing value would influence the analysis result by acting as an improper feature that mislead the modeling process greatly and thus lead to a incorrect result. Therefore, we need to take care the missing value properly and try to fill out the missing value with proper and reasonable value.

- Single imputation method
- Regression imputation
- Model based imputation

Correlation



Method & Model



Method & Model



ExtraTreesRegressor Model

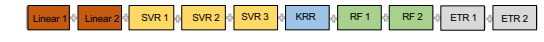
The cross validation score is 0.858874576346.

This model uses ensemble method which can combine the different models and reduce the bias.

The figure shows the predicted values are close to real values.

Model Average

Model Averaging



10

Achievement

305	30	Zheko	0.11623	3	10d
306	259	Zaur	0.11623	23	7d
307	new	Data Mining Project	0.11623	1	now

Your Best Entry Your submission scored 0.11623, which is not an improvement of your best score. Keep trying!

Summary

Through this project, we realized that a real-world data set which involves complex features is much more complicated than some general data sets like MNIST. Prior to performing regression training, there are lots of feature engineering jobs to do which definitely have large impact on our prediction model. Also the model selection is not pretty deterministic. We see from the experiments that lasso regression could easily get good performance but other approaches like random forest need to well tune hyper-parameters to get decent results.