

Machine Learning Engineer Nanodegree

Capstone Project-Kaggle House Prices: Advanced Regression Techniques

Süleyman Diker

July 07st, 2018

I. Definition

Estimates will be made about the monetary value of a house over the most suitable model using real house property and price data. After an exhaustive feature engineering, multiple machine learning regression algorithms will be trained and tuned, to obtain a better predictive performance.

Project Overview

The purpose of this project is to estimate the price of house in Iowa with the help of data containing 81 property definitions. The data for this project was obtained from Kaggle in the House Prices: Advanced Regression Techniques competition.

Problem Statement

Our main goal is to estimate the house price, which we have with the data we have, as accurate as possible. We can list the most basic items that we need to solve in the framework of this problem as follows:

- Exploratory data analysis: The nature of the features in the data is crucial for accurate estimation of missing values of the existing properties, number of outliers and detection of outliers.
- Feature preprocessing: It is a process of visualizing the data to detect if there is an inconsistency or incomplete data, and to process the data in this direction.
- Benchmark modeling: To solve the problem we start by creating a model using a standard technique.
- Model improvement: We try to optimize our model by setting parameters and using learning techniques to improve the model's performance

Metrics

In this project we will measure the success of our model using the RMSLE (Root Mean Squared Logarithmic Error) metrics, which calculate the difference between the actual price and the model predicted by our model.

As the difference between the actual value of the house and the estimated value increases, the value of RMSLE also increases, so we need to have the minimum RMSLE value.

II. Analysis

Data Exploration

The training set includes categorical and numerical features as well as some missing values. Since we often can not have a complete dataset in every direction, it is important that we do a good analysis of the dataset. Our data set is divided into two parts as training set and test set. The training set consists of 1460 rows with 81

features. The test set consists of 1459 rows with 80 features. Our test dataset is missing a column according to our training dataset. For this reason, we will guess the SalePrice column by modeling the data through our training dataset.

When we examine the house price histogram in Figure 1, we see that the skewed value is high. (Skew Value: 1.8828757597682129) For this reason, we need log transformations to normalize the skewed value, so we use RMSLE as our evaluation metric. After the log transformations, we see that the skewed value falls and that the Figure 2 is a more balanced graph. (Skew Value: 0.12133506220520406)

In [1]:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 train = pd.read_csv('train.csv')
6 test = pd.read_csv('test.csv')
7
8 print ("Train data shape:", train.shape)
9 print ("Test data shape:", test.shape)
```

```
('Train data shape:', (1460, 81))
('Test data shape:', (1459, 80))
```

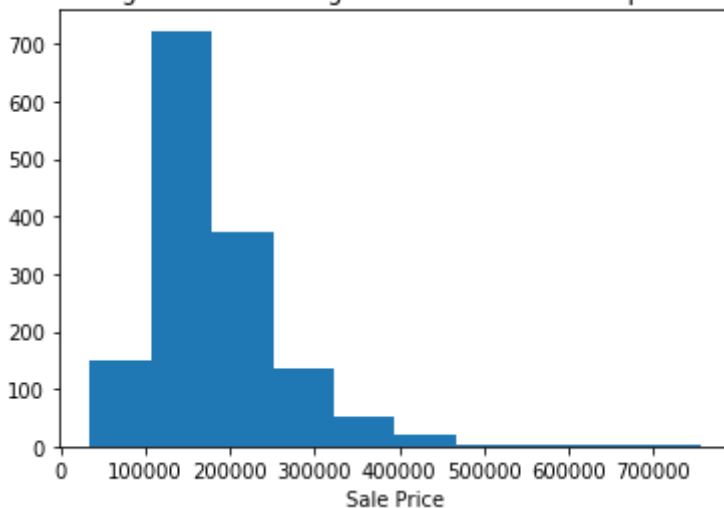
In [2]:

```

1
2 plt.xlabel("Sale Price")
3 plt.title("Figure 1: The histograms of the house sale price")
4 plt.hist(train.SalePrice)
5 plt.show()
6 print ("Sale Price skew is:", train.SalePrice.skew())
7
8 target = np.log(train.SalePrice)
9 plt.title("Figure 2: Original data; (right) Log transformed data.")
10 plt.xlabel("Log(Sale Price)")
11 plt.hist(target)
12 plt.show()
13 print ("Log(Sale Price) skew is:", target.skew())

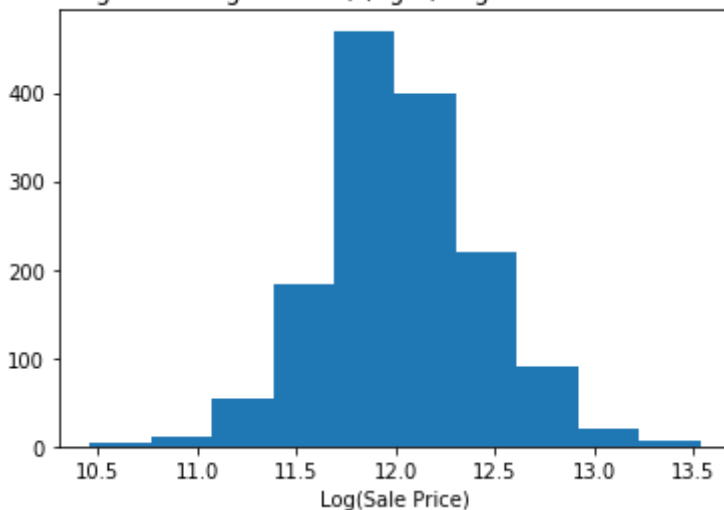
```

Figure 1: The histograms of the house sale price



('Sale Price skew is:', 1.8828757597682129)

Figure 2: Original data; (right) Log transformed data.



('Log(Sale Price) skew is:', 0.12133506220520406)

Exploratory Visualization

The correlation values between the numerical variables and the home sales are as follows. We can not say that each of the feature has the same effect on the sales value of the house, each feature has positive, negative or neutral effects on the sales value.

In [3]:

```

1 numeric_features = train.select_dtypes(include=[np.number])
2 numeric_features.dtypes
3 corr = numeric_features.corr()
4
5 print (corr['SalePrice'].sort_values(ascending=False)[0:], '\n')

```

```

(SalePrice      1.000000
OverallQual    0.790982
GrLivArea      0.708624
GarageCars     0.640409
GarageArea     0.623431
TotalBsmntSF   0.613581
1stFlrSF       0.605852
FullBath       0.560664
TotRmsAbvGrd   0.533723
YearBuilt      0.522897
YearRemodAdd    0.507101
GarageYrBlt    0.486362
MasVnrArea     0.477493
Fireplaces     0.466929
BsmntFinSF1    0.386420
LotFrontage    0.351799
WoodDeckSF     0.324413
2ndFlrSF       0.319334
OpenPorchSF    0.315856
HalfBath       0.284108
LotArea        0.263843
BsmntFullBath  0.227122
BsmntUnfSF     0.214479
BedroomAbvGr   0.168213
ScreenPorch    0.111447
PoolArea       0.092404
MoSold         0.046432
3SsnPorch      0.044584
BsmntFinSF2    -0.011378
BsmntHalfBath  -0.016844
MiscVal        -0.021190
Id             -0.021917
LowQualFinSF   -0.025606
YrSold         -0.028923
OverallCond    -0.077856
MSSubClass     -0.084284
EnclosedPorch  -0.128578
KitchenAbvGr   -0.135907
Name: SalePrice, dtype: float64, '\n')

```

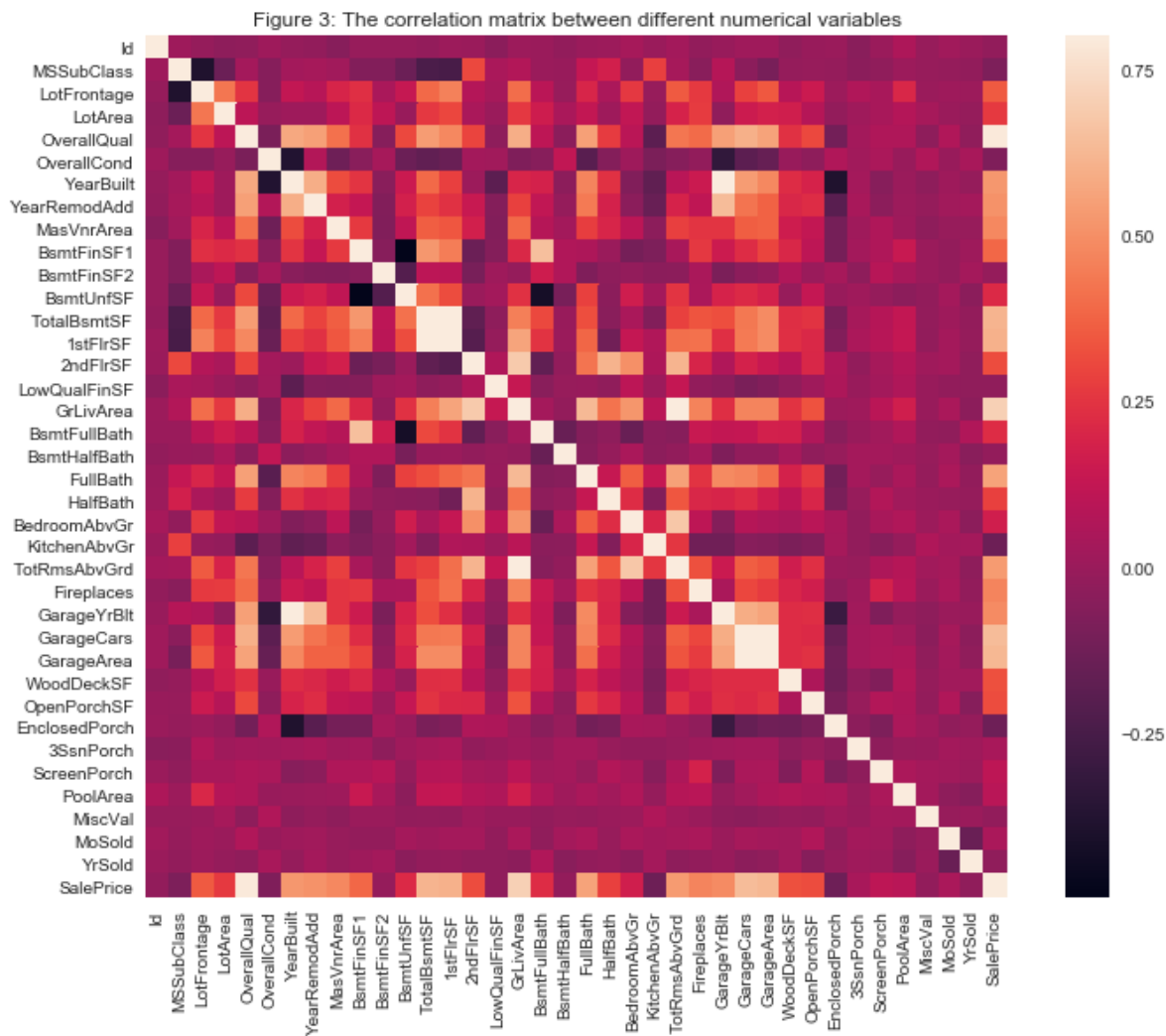
We can see the correlation matrix between the different numerical variables below:

In [4]:

```

1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 corr_matrix = train.corr()
5 plt.subplots
6 sns.set(rc={'axes.facecolor':'white', 'figure.facecolor':'white'})
7 f, ax = plt.subplots(figsize=(12, 9))
8 sns.heatmap(corr_matrix, vmax=.8, square=True)
9 sns.heatmap
10 plt.title("Figure 3: The correlation matrix between different numerical variables")
11 plt.show()

```



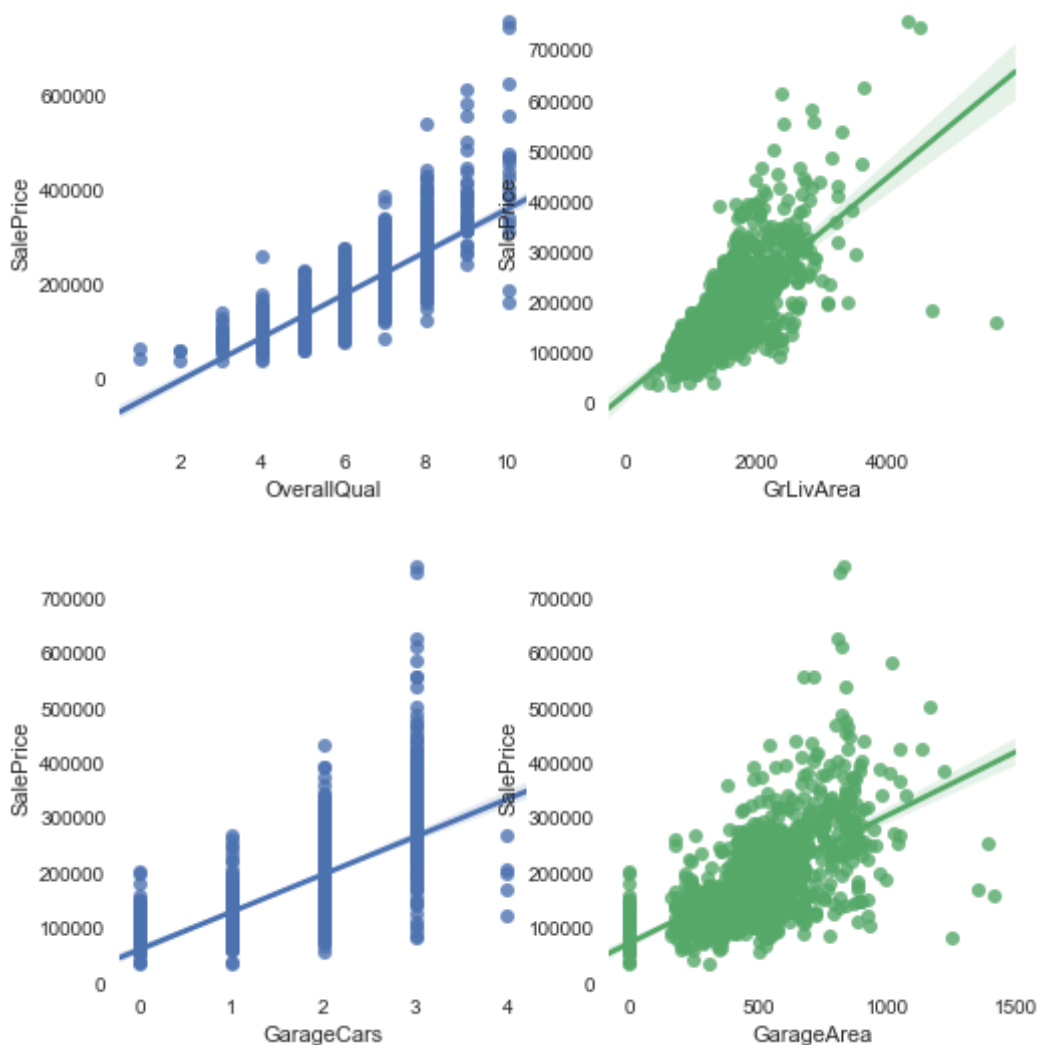
We can look at the top 6 feature distribution charts that have the strongest correlation that affects home selling prices:

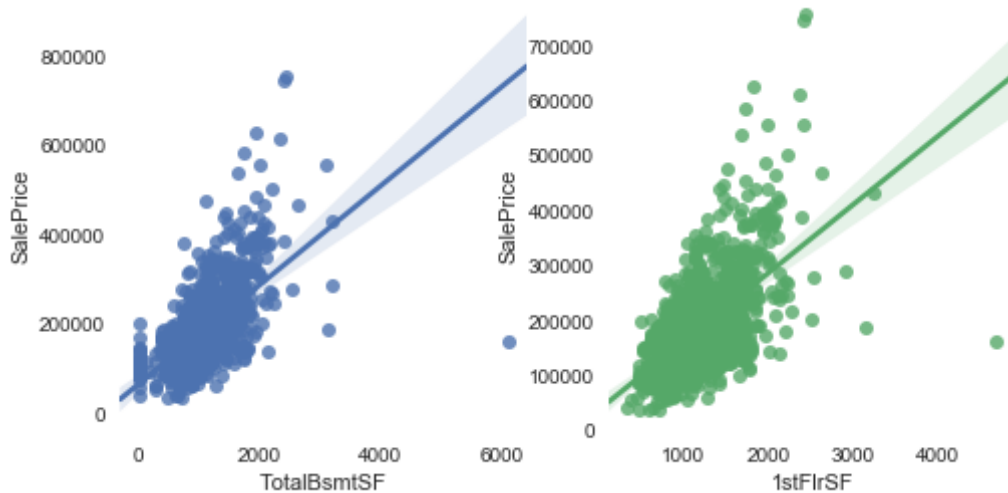
In [5]:

```

1 fig1, (ax1, ax2) = plt.subplots(ncols=2,figsize = (8, 4))
2 sns.regplot(x = 'OverallQual', y = 'SalePrice', data = train,ax=ax1)
3 sns.regplot(x = 'GrLivArea', y = 'SalePrice', data = train,ax=ax2)
4 plt.show()
5
6 fig2, (ax1, ax2) = plt.subplots(ncols=2,figsize = (8, 4))
7 sns.regplot(x = 'GarageCars', y = 'SalePrice', data = train,ax=ax1)
8 sns.regplot(x = 'GarageArea', y = 'SalePrice', data = train,ax=ax2)
9 plt.show()
10
11 fig3, (ax1, ax2) = plt.subplots(ncols=2,figsize = (8, 4))
12 sns.regplot(x = 'TotalBsmtSF', y = 'SalePrice', data = train,ax=ax1)
13 sns.regplot(x = '1stFlrSF', y = 'SalePrice', data = train,ax=ax2)
14 plt.show()
15
16

```



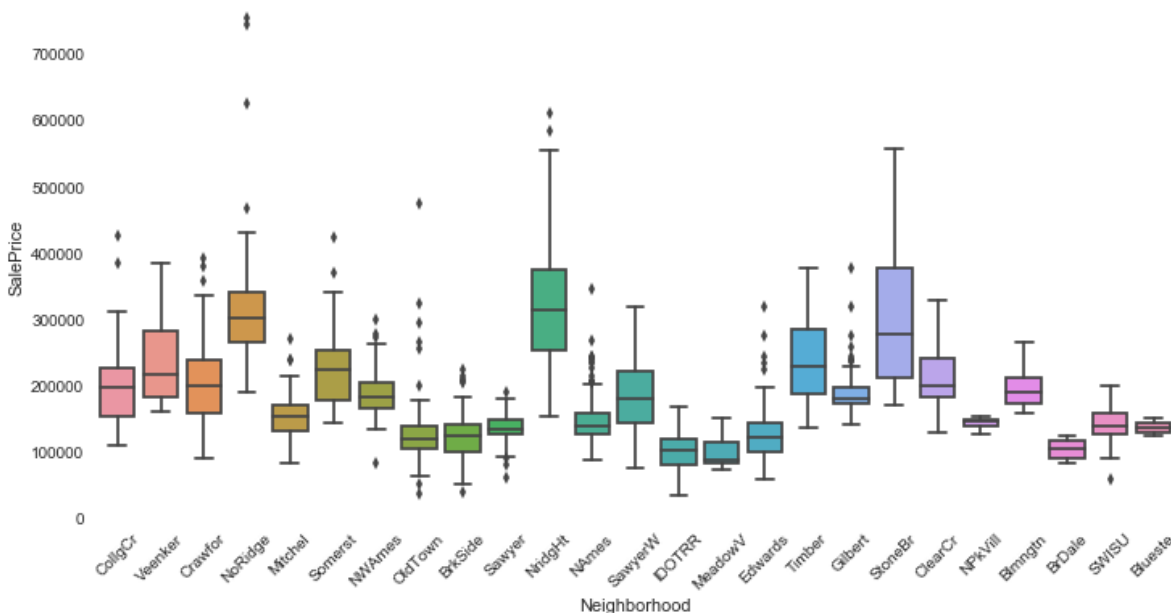


The relationship between the neighborhood and the resale price is clearly seen in Figure 5..We observe that the neighborhood feature is a feature that seriously affects the price of the house. For this reason, in categorical variables, all values should be evaluated as equal.If values such as 0 and 1 are given to convert to number, we will not reflect the effect of this variable correctly on the model.So we should pay particular attention to the transformations of categorical variables.

In [6]:

```
1 plt.figure(figsize = (12, 6))
2 sns.boxplot(x = 'Neighborhood', y = 'SalePrice', data = train)
3 xt = plt.xticks(rotation=45)
4 plt.title("Figure 4:Boxplot of sale price in different neighborhood.")
5 plt.show()
```

Figure 4:Boxplot of sale price in different neighborhood.



Algorithms and Techniques

Our main goal is to make the best guess scoring by doing experiments within the scope of advanced regression techniques, so I tried the community learning method.This technique aims to learn the weighted combinations of basic models.

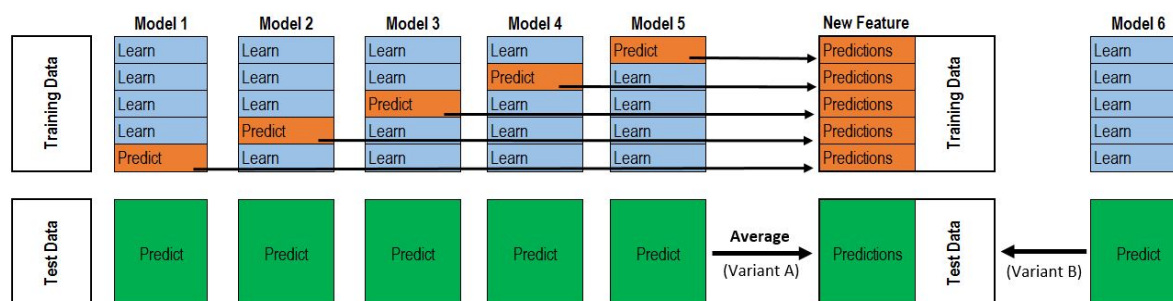
This technique aims to learn the weighted combinations of basic models. This technique is also called technical stacking. This technique combines all model predictions to minimize the generalization error and tries to capture the most accurate prediction score.

we see two levels of five-fold stacking schemes as schematics below. First, the training dataset need to be split into 5 folds. Second, we iterate over this 5 folds training dataset. In each iteration, each base model will be trained using 4 folds and predict on the hold out fold. At the same time, each base model also need to provide a prediction on the entire test dataset. After the iteration over all folds, we will have the prediction of the entire training dataset for each model and 5 copies of the prediction of the entire test dataset for each model. Finally, we train second level model, or stacker, using the prediction in the training dataset as new features and use the average of the 5 copies of the test dataset predictions as the test input for the trained model to provide the final prediction.

In [7]:

```
1 from IPython.display import Image
2 Image(filename='images/QBuDOjs.jpg')
```

Out[7]:



Benchmark

The Random Forest regression is selected as our benchmark model for this project. The default parameters in scikit-learn package will be used as the benchmark.

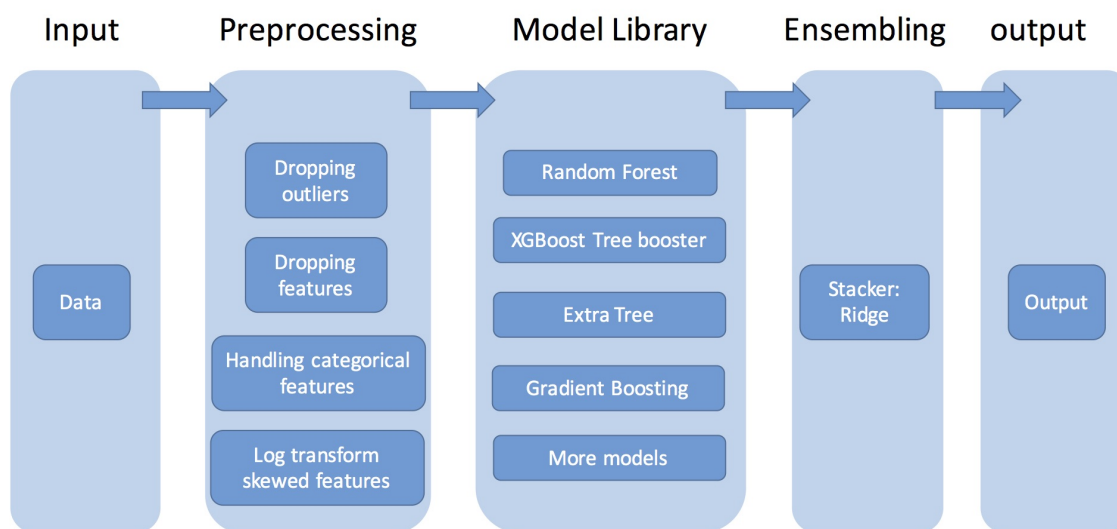
III. Methodology

In general, the methodology and modeling approach is illustrated in the flow chart below.

In [8]:

```
1 from IPython.display import Image
2 Image(filename='images/flowchart.jpg')
```

Out[8]:



Data Preprocessing

We can summarize the data pre-processing operations as follows:

- Outliers deletion: We examined the data intensities of the features with the highest correlation effect. We have found data that are very different from the data intensities of these properties and we have eliminated the data as follows:

In [21]:

```

1 outlier_idx = []
2
3 for i in train[train['GrLivArea'] > 4000]['GrLivArea'].index:
4     outlier_idx.append(i)
5
6 for i in train[train['GarageArea'] > 1200]['GarageArea'].index:
7     outlier_idx.append(i)
8
9
10 for i in train[train['TotalBsmtSF'] > 3000]['TotalBsmtSF'].index:
11     outlier_idx.append(i)
12
13 for i in train[train['1stFlrSF'] > 3000]['1stFlrSF'].index:
14     outlier_idx.append(i)
15
16 for i in train[train['MasVnrArea'] > 800]['MasVnrArea'].index:
17     outlier_idx.append(i)
18
19 for i in train[train['BsmtFinSF1'] > 2000]['BsmtFinSF1'].index:
20     outlier_idx.append(i)
21
22 for i in train[train['LotFrontage'] > 150]['LotFrontage'].index:
23     outlier_idx.append(i)
24
25 for i in train[train['WoodDeckSF'] > 600]['WoodDeckSF'].index:
26     outlier_idx.append(i)
27
28 for i in train[train['2ndFlrSF'] > 1650]['2ndFlrSF'].index:
29     outlier_idx.append(i)
30
31 for i in train[train['OpenPorchSF'] > 400]['OpenPorchSF'].index:
32     outlier_idx.append(i)
33
34 for i in train[train['LotArea'] > 50000]['LotArea'].index:
35     outlier_idx.append(i)
36
37 for i in train[train['BsmtUnfSF'] > 2000]['BsmtUnfSF'].index:
38     outlier_idx.append(i)
39
40 print(outlier_idx)
41
42 train.drop(train.index[outlier_idx],inplace=True)
43
44

```

```

[523, 691, 1182, 1298, 581, 1298, 440, 523, 1298, 523, 1298, 70, 224,
 349, 691, 798, 1169, 523, 898, 1182, 1298, 1182, 1298, 691, 1169, 118
2, 523, 313, 1298, 224, 477, 581]

```

- Feature selection: The categorical variables with the greatest number of missing values are Alley, FireplaceQu, PoolQC, Fence and MiscFeature. We can see this clearly in the following table. We delete these features with very high amounts of missing data as follows:

In [10]:

```

1 total = train.isnull().sum().sort_values(ascending = False)
2 percent = (train.isnull().sum() / train.isnull().count()).sort_values(ascending
3 missing_data = pd.concat([total, percent], axis=1, keys=['Total', 'Percent'])
4 missing_data.head(10)

```

Out[10]:

	Total	Percent
PoolQC	1402	0.997155
MiscFeature	1357	0.965149
Alley	1320	0.938834
Fence	1132	0.805121
FireplaceQu	684	0.486486
LotFrontage	251	0.178521
GarageCond	79	0.056188
GarageType	79	0.056188
GarageYrBlt	79	0.056188
GarageFinish	79	0.056188

In [11]:

```

1 all_data = pd.concat((train.loc[:, 'MSSubClass': 'SaleCondition'],
2                        test.loc[:, 'MSSubClass': 'SaleCondition']))
3
4 to_delete = ['Alley', 'FireplaceQu', 'PoolQC', 'Fence', 'MiscFeature']
5 all_data = all_data.drop(to_delete, axis=1)

```

- Log transform skewed features: All features that have a skewness larger than 0.75 are log transformed.

In [13]:

```

1 from scipy.stats import skew
2
3 numeric_feats = all_data.dtypes[all_data.dtypes != "object"].index
4 skewed_feats = train[numeric_feats].apply(lambda x: skew(x.dropna()))
5 skewed_feats = skewed_feats[skewed_feats > 0.75]
6 skewed_feats = skewed_feats.index
7 all_data[skewed_feats] = np.log1p(all_data[skewed_feats])

```

- Handle categorical features: The categorical features are handled by the pandas get_dummies function that converts categorical variables into dummy or indicator variables. We do this transformation in the following way:

In [14]:

```

1 all_data = pd.get_dummies(all_data)
2 all_data = all_data.fillna(all_data.mean())

```

Implementation

Before applying the ensemble learning method, I trained and tuned 4 base model including Random Forest regressor, Extra Trees regressor, Gradient Boosting regressor, Extreme Gradient Boosting regressor.

At the beginning of the ensembling implementation, I use only 4 base models to form my model library, including Random Forest regressor, Extra Trees regressor, Gradient Boosting regressor, Extreme Gradient Boosting regressor. And I use the Ridge Regression as my second level model. The idea of the stacking approach is to use the predictions from the first level model as the new features for the second level model and make prediction based on that.

The base models should be as unrelated as possible and this is why we tend to include more models in the ensemble even though they may not perform well. For my final, I use a much larger much library including Random Forest regressor, Extra Trees regressor, Gradient Boosting regressor, Extreme Gradient Boosting regressor, LassoCV regressor, K neighbors regressor, LassoLarsCV regressor, Elastic Net regressor, Support Vector Machine regressor. The Extreme Gradient Boosting regressor is from the XGBoost package, all other techniques are from scikit-learn package. The basic idea of the final model is to use a larger base model library in order to achieve better results. Some of the base models are used only with their scikit-learn default parameter values without fine tuning.

There are two challenging parts during the coding process. The first one is that we need to make sure the feature for the final second level prediction is based on the prediction of the first level base models. Thus, we also need to use the model prediction as the features for the test dataset. The second one is how exactly should we choose the base models. This is a tricky and not fully addressed question in this project since there are so many models to select.

Refinement

The refinement is conducted using the grid search technique. For the learning purpose of this project, I only use a few tunable parameters. Before the parameter tuning, I just use the scikit-learn default parameter values for each model. After the refinement of the base models, we using stacking technique to ensembling them together to provide the final submission.

Random forest regression: Best Params:{} Best CV Score: 0.146323732016

eXtreme Gradient Boosting regression: Best Params:{} Best CV Score: 0.127444960531

Extra trees regression: Best Params:{} Best CV Score: 0.149811744332

Gradient boosted tree regression: Best Params:{} Best CV Score: 0.126801721411

IV. Results

Model Evaluation and Validation

One interesting observation during the model ensembling process is that increasing the size of the model library can potentially increase both the validation score and the real score. A variety of base models are included in the model library with their scikit-learn default parameter values. I need to acknowledge that the approach is somewhat brute force and a more elegant way to do this is to find the base models to be as uncorrelated as possible.

Justification

In summary, the ensemble approach in this project greatly improve the performance of the model in terms of the Kaggle Leaderboard score and rank from the benchmark random forecast regression ranked. The result is quite encouraging and opens a lot more to think, learn and explore.

V. Conclusion

Free-Form Visualization

As my ultimate goal in this project is to build a model as accurate as possible measure by the Kaggle Leaderboard score, I will use the Kaggle Leaderboard score at the different model development stage as my final project visualization. This project opens several avenues for future work. One part is for data preprocessing, such as more advanced outlier detection techniques, more advanced feature selection techniques, fine tuning of the log transform threshold of the skewed variable and so on. More creative feature engineering will be also extremely useful. Another part is for the regression techniques. My intuition for ensemble learning technique is that a better score may be achieved by increasing the size of the model library. The dilemma here is a tradeoff between the amount of time and engineering efforts you want to spend and the actual value of the tiny improvement.

Reflection

For this project, I started with an exploratory data analysis, in which I found that there are some features being skewed, some features are categorical, and some missing values. Several data preprocessing steps were conducted in order to make the data ready for the machine learning modeling part. For the modeling part, I firstly trained several base models and tuned their parameters. Then, I applied the ensemble learning method to stack all the base models together to provide the final submission. One of difficult parts is data preprocessing since it is hard to explain in a rigorous way of why I did a specific choice to improve the model. There are a lot of exploratory tests and investigations to justify the choice only based on the local cross validation scores. The other challenge is the stacking approach itself since it is more complicated than a single model approach. It is quite easy to use the wrong features for the stacker since it makes prediction based on the predictions of all the base models. The key to success is to write a pseudocode and check it carefully before writing the production code.

Improvement

This project opens several avenues for future work. One part is for data preprocessing, such as more advanced outlier detection techniques, more advanced feature selection techniques, fine tuning of the log transform threshold of the skewed variable and so on. More creative feature engineering will be also extremely useful. Another part is for the regression techniques. My intuition for ensemble learning technique is that a better score may be achieved by increasing the size of the model library. The dilemma here is a tradeoff between the amount of time and engineering efforts you want to spend and the actual value of the tiny improvement.

In [23]:

```
1 import numpy as np
2 import pandas as pd
3 import datetime
4 from sklearn.cross_validation import KFold
5 from sklearn.cross_validation import train_test_split
6 import time
7 from sklearn import preprocessing
8 from xgboost import XGBRegressor
9 from sklearn.ensemble import RandomForestRegressor, ExtraTreesRegressor, GradientBoostingRegressor
10 from sklearn.grid_search import GridSearchCV
11 from sklearn.cross_validation import ShuffleSplit
12 from sklearn.metrics import make_scorer, mean_squared_error
13 from sklearn.linear_model import Ridge, LassoCV, LassoLarsCV, ElasticNet
14 from sklearn.kernel_ridge import KernelRidge
15 from sklearn.neighbors import KNeighborsRegressor
16 from sklearn.svm import SVR
17 from scipy.stats import skew
18
19 def create_submission(prediction, score):
20     now = datetime.datetime.now()
21     sub_file = 'submission_'+str(score)+'_'+str(now.strftime("%Y-%m-%d-%H-%M"))
22     print ('Creating submission: ', sub_file)
23     pd.DataFrame({'Id': test['Id'].values, 'SalePrice': prediction}).to_csv(sub_file+'.csv')
24
25
26 def data_preprocess(train, test):
27
28     outlier_idx = []
29
30     for i in train[train['GrLivArea'] > 4000]['GrLivArea'].index:
31         outlier_idx.append(i)
32
33     for i in train[train['GarageArea'] > 1200]['GarageArea'].index:
34         outlier_idx.append(i)
35
36
37     for i in train[train['TotalBsmtSF'] > 3000]['TotalBsmtSF'].index:
38         outlier_idx.append(i)
39
40     for i in train[train['1stFlrSF'] > 3000]['1stFlrSF'].index:
41         outlier_idx.append(i)
42
43     for i in train[train['MasVnrArea'] > 800]['MasVnrArea'].index:
44         outlier_idx.append(i)
45
46     for i in train[train['BsmtFinSF1'] > 2000]['BsmtFinSF1'].index:
47         outlier_idx.append(i)
48
49     for i in train[train['LotFrontage'] > 150]['LotFrontage'].index:
50         outlier_idx.append(i)
51
52     for i in train[train['WoodDeckSF'] > 600]['WoodDeckSF'].index:
53         outlier_idx.append(i)
54
55     for i in train[train['2ndFlrSF'] > 1650]['2ndFlrSF'].index:
56         outlier_idx.append(i)
57
58     for i in train[train['OpenPorchSF'] > 400]['OpenPorchSF'].index:
59         outlier_idx.append(i)
```

```

60
61     for i in train[train['LotArea'] > 50000]['LotArea'].index:
62         outlier_idx.append(i)
63
64     for i in train[train['BsmtUnfSF'] > 2000]['BsmtUnfSF'].index:
65         outlier_idx.append(i)
66
67     print(outlier_idx)
68
69     train.drop(train.index[outlier_idx],inplace=True)
70     all_data = pd.concat((train.loc[:,'MSSubClass':'SaleCondition'],
71                          test.loc[:,'MSSubClass':'SaleCondition']))
72
73     to_delete = ['Alley','FireplaceQu','PoolQC','Fence','MiscFeature']
74     all_data = all_data.drop(to_delete,axis=1)
75     train["SalePrice"] = np.log1p(train["SalePrice"])
76
77     numeric_feats = all_data.dtypes[all_data.dtypes != "object"].index
78     skewed_feats = train[numeric_feats].apply(lambda x: skew(x.dropna())) #comp
79     skewed_feats = skewed_feats[skewed_feats > 0.75]
80     skewed_feats = skewed_feats.index
81     all_data[skewed_feats] = np.log1p(all_data[skewed_feats])
82     all_data = pd.get_dummies(all_data)
83     all_data = all_data.fillna(all_data.mean())
84     X_train = all_data[:train.shape[0]]
85     X_test = all_data[train.shape[0]:]
86     y = train.SalePrice
87
88     return X_train,X_test,y
89
90
91 def mean_squared_error_(ground_truth, predictions):
92     return mean_squared_error(ground_truth, predictions) ** 0.5
93 RMSE = make_scorer(mean_squared_error_, greater_is_better=False)
94
95 class ensemble(object):
96     def __init__(self, n_folds, stacker, base_models):
97         self.n_folds = n_folds
98         self.stacker = stacker
99         self.base_models = base_models
100     def fit_predict(self,train,test,ytr):
101         X = train.values
102         y = ytr.values
103         T = test.values
104         folds = list(KFold(len(y), n_folds = self.n_folds, shuffle = True, rand
105         S_train = np.zeros((X.shape[0],len(self.base_models)))
106         S_test = np.zeros((T.shape[0],len(self.base_models)))
107         for i,reg in enumerate(base_models):
108             print ("Fitting the base model...")
109             S_test_i = np.zeros((T.shape[0],len(folds)))
110             for j, (train_idx,test_idx) in enumerate(folds):
111                 X_train = X[train_idx]
112                 y_train = y[train_idx]
113                 X_holdout = X[test_idx]
114                 reg.fit(X_train,y_train)
115                 y_pred = reg.predict(X_holdout)[:]
116                 S_train[test_idx,i] = y_pred
117                 S_test_i[:,j] = reg.predict(T)[:]
118             S_test[:,i] = S_test_i.mean(1)
119
120

```

```

121     print ("Stacking base models...")
122     # tuning the stacker
123     param_grid = {
124         'alpha': [1e-3,5e-3,1e-2,5e-2,1e-1,0.2,0.3,0.4,0.5,0.8,1e0,3,5,7,1e1],
125     }
126     grid = GridSearchCV(estimator=self.stacker, param_grid=param_grid, n_jobs=1)
127     grid.fit(S_train, y)
128     try:
129         print('Param grid:')
130         print(param_grid)
131         print('Best Params:')
132         print(grid.best_params_)
133         print('Best CV Score:')
134         print(-grid.best_score_)
135         print('Best estimator:')
136         print(grid.best_estimator_)
137         print(message)
138     except:
139         pass
140
141     y_pred = grid.predict(S_test)[: ]
142     return y_pred, -grid.best_score_
143
144 train = pd.read_csv("train.csv") # read train data
145 test = pd.read_csv("test.csv") # read test data
146
147 # build a model library (can be improved)
148 base_models = [
149     RandomForestRegressor(
150         n_jobs=1, random_state=0,
151         n_estimators=500, max_features=14
152     ),
153     RandomForestRegressor(
154         n_jobs=1, random_state=0,
155         n_estimators=500, max_features=20,
156     max_depth = 7
157     ),
158     ExtraTreesRegressor(
159         n_jobs=1, random_state=0,
160         n_estimators=500, max_features=15
161     ),
162     ExtraTreesRegressor(
163         n_jobs=1, random_state=0,
164         n_estimators=500, max_features=20
165     ),
166     GradientBoostingRegressor(
167         random_state=0,
168         n_estimators=500, max_features=10, max_depth=6,
169         learning_rate=0.05, subsample=0.8
170     ),
171     GradientBoostingRegressor(
172         random_state=0,
173         n_estimators=500, max_features=15, max_depth=6,
174         learning_rate=0.05, subsample=0.8
175     ),
176     XGBRegressor(
177         seed=0,
178         n_estimators=500, max_depth=10,
179         learning_rate=0.05, subsample=0.8, colsample_bytree=0.75
180     ),
181

```



```

182         XGBRegressor(
183             seed=0,
184             n_estimators=500, max_depth=7,
185             learning_rate=0.05, subsample=0.8, colsample_bytree=0.75
186         ),
187         LassoCV(alphas = [1, 0.1, 0.001, 0.0005]),
188         KNeighborsRegressor(n_neighbors = 5),
189         KNeighborsRegressor(n_neighbors = 10),
190         KNeighborsRegressor(n_neighbors = 15),
191         KNeighborsRegressor(n_neighbors = 25),
192         LassoLarsCV(),
193         ElasticNet(),
194         SVR()
195     ]
196
197     ensem = ensemble(
198         n_folds=5,
199         stacker=Ridge(),
200         base_models=base_models
201     )
202
203     X_train,X_test,y_train = data_preprocess(train,test)
204     y_pred, score = ensem.fit_predict(X_train,X_test,y_train)
205
206     create_submission(np.expml(y_pred),score)

```

```

/Users/tcsdiker/anaconda2/lib/python2.7/site-packages/sklearn/linear_model/least_angle.py:313: ConvergenceWarning: Regressors in active set degenerate. Dropping a regressor, after 112 iterations, i.e. alpha=1.266e-04, with an active set of 102 regressors, and the smallest cholesky pivot element being 1.490e-08. Reduce max_iter or increase eps parameters.

```

```
ConvergenceWarning)
```

```

/Users/tcsdiker/anaconda2/lib/python2.7/site-packages/sklearn/linear_model/least_angle.py:313: ConvergenceWarning: Regressors in active set degenerate. Dropping a regressor, after 132 iterations, i.e. alpha=8.659e-05, with an active set of 122 regressors, and the smallest cholesky pivot element being 1.490e-08. Reduce max_iter or increase eps parameters.

```

```
ConvergenceWarning)
```

```

/Users/tcsdiker/anaconda2/lib/python2.7/site-packages/sklearn/linear_model/least_angle.py:313: ConvergenceWarning: Regressors in active set degenerate. Dropping a regressor, after 132 iterations, i.e. alpha=8.659e-05, with an active set of 122 regressors, and the smallest cholesky pivot element being 1.054e-08. Reduce max_iter or increase eps parameters.

```

In [20]:

```

1 import datetime
2 import numpy as np
3 import pandas as pd
4 import xgboost as xgb
5 from sklearn import preprocessing
6 from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor,
7 from sklearn.grid_search import GridSearchCV
8 from sklearn.cross_validation import ShuffleSplit
9 from sklearn.metrics import make_scorer, mean_squared_error
10 from sklearn.kernel_ridge import KernelRidge
11 from sklearn.svm import SVR
12 from sklearn.neighbors import KNeighborsRegressor
13 from scipy.stats import skew
14
15 def mean_squared_error_(ground_truth, predictions):
16     return mean_squared_error(ground_truth, predictions) ** 0.5
17 RMSE = make_scorer(mean_squared_error_, greater_is_better=False)
18
19 def create_submission(prediction,score):
20     now = datetime.datetime.now()
21     sub_file = 'submission_'+str(score)+'_'+str(now.strftime("%Y-%m-%d-%H-%M"))
22     print ('Creating submission: ', sub_file)
23     pd.DataFrame({'Id': test['Id'].values, 'SalePrice': prediction}).to_csv(sub
24
25 def data_preprocess(train,test):
26     train.drop(train.index[outlier_idx],inplace=True)
27     all_data = pd.concat((train.loc[:, 'MSSubClass': 'SaleCondition'],
28                           test.loc[:, 'MSSubClass': 'SaleCondition']))
29
30     to_delete = ['Alley', 'FireplaceQu', 'PoolQC', 'Fence', 'MiscFeature']
31     all_data = all_data.drop(to_delete,axis=1)
32
33     train["SalePrice"] = np.log1p(train["SalePrice"])
34     #log transform skewed numeric features
35     numeric_feats = all_data.dtypes[all_data.dtypes != "object"].index
36     skewed_feats = train[numeric_feats].apply(lambda x: skew(x.dropna())) #comp
37     skewed_feats = skewed_feats[skewed_feats > 0.75]
38     skewed_feats = skewed_feats.index
39     all_data[skewed_feats] = np.log1p(all_data[skewed_feats])
40     all_data = pd.get_dummies(all_data)
41     all_data = all_data.fillna(all_data.mean())
42     X_train = all_data[:train.shape[0]]
43     X_test = all_data[train.shape[0]:]
44     y = train.SalePrice
45
46     return X_train,X_test,y
47
48 def model_random_forecast(Xtrain,Xtest,ytrain):
49
50     X_train = Xtrain
51     y_train = ytrain
52     rfr = RandomForestRegressor(n_jobs=1, random_state=0)
53     param_grid = {}
54     model = GridSearchCV(estimator=rfr, param_grid=param_grid, n_jobs=1, cv=10,
55 model.fit(X_train, y_train)
56     print('Random forecast regression...')
57     print('Best Params:')
58     print(model.best_params_)
59     print('Best CV Score:')

```

```
60     print(-model.best_score_)
61
62     y_pred = model.predict(Xtest)
63     return y_pred, -model.best_score_
64
65 def model_gradient_boosting_tree(Xtrain,Xtest,ytrain):
66
67     X_train = Xtrain
68     y_train = ytrain
69     gbr = GradientBoostingRegressor(random_state=0)
70     param_grid = {}
71     model = GridSearchCV(estimator=gbr, param_grid=param_grid, n_jobs=1, cv=10,
72     model.fit(X_train, y_train)
73     print('Gradient boosted tree regression...')
74     print('Best Params:')
75     print(model.best_params_)
76     print('Best CV Score:')
77     print(-model.best_score_)
78
79     y_pred = model.predict(Xtest)
80     return y_pred, -model.best_score_
81
82 def model_xgb_regression(Xtrain,Xtest,ytrain):
83
84     X_train = Xtrain
85     y_train = ytrain
86
87     xgbreg = xgb.XGBRegressor(seed=0)
88     param_grid = {}
89     model = GridSearchCV(estimator=xgbreg, param_grid=param_grid, n_jobs=1, cv=
90     model.fit(X_train, y_train)
91     print('eXtreme Gradient Boosting regression...')
92     print('Best Params:')
93     print(model.best_params_)
94     print('Best CV Score:')
95     print(-model.best_score_)
96
97     y_pred = model.predict(Xtest)
98     return y_pred, -model.best_score_
99
100 def model_extra_trees_regression(Xtrain,Xtest,ytrain):
101
102     X_train = Xtrain
103     y_train = ytrain
104
105     etr = ExtraTreesRegressor(n_jobs=1, random_state=0)
106     param_grid = {}
107     model = GridSearchCV(estimator=etr, param_grid=param_grid, n_jobs=1, cv=10,
108     model.fit(X_train, y_train)
109     print('Extra trees regression...')
110     print('Best Params:')
111     print(model.best_params_)
112     print('Best CV Score:')
113     print(-model.best_score_)
114
115     y_pred = model.predict(Xtest)
116     return y_pred, -model.best_score_
117
118
119 # read data, build model and do prediction
120 train = pd.read_csv("train.csv") # read train data
```