

141C Final Project

Team members:

Deng Zeyu
Ricardo Rendon
Julia Stelman
Wei-Kaung Lin

Project topic:

House price prediction challenge

<https://www.kaggle.com/c/house-prices-advanced-regression-techniques>

I. Data Exploration

In the dataset, there are 19 discrete variables. There are 27 categorical variables which we have formatted as ordinal. There are 32 non-ordinal categorical variables. There are 3 continuous variables.

Missing value

Missing value percentage of variables: MiscFeature have more than 50% of missing values, so it was removed from prediction. For most other variables, a mean value was used to replace the missing value. The exception was GarageYrBlt, for which missing values were replaced with the mean - 3.5* standard deviation.

	Variable	Missing Percentage
3	LotFrontage	0.177397
25	MasVnrType	0.005479
26	MasVnrArea	0.005479
42	Electrical	0.000685
59	GarageYrBlt	0.055479
74	MiscFeature	0.963014

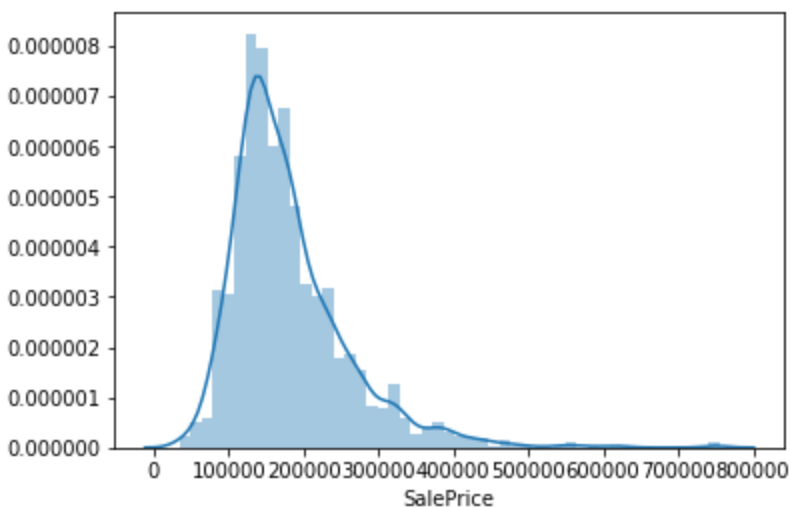
Analysing distribution of the response variable, Sales Price:

First, we want to analyse the sales price in this area. Housing price is the most important thing we need to consider when people buy a house.

Here is our result:

```
count    1460.000000
mean     180921.195890
std      79442.502883
min       34900.000000
25%      129975.000000
50%      163000.000000
75%      214000.000000
max       755000.000000
```

We have found the minimum price is 34900, it is larger than 0, so we can say there is no impossible data on this variable.

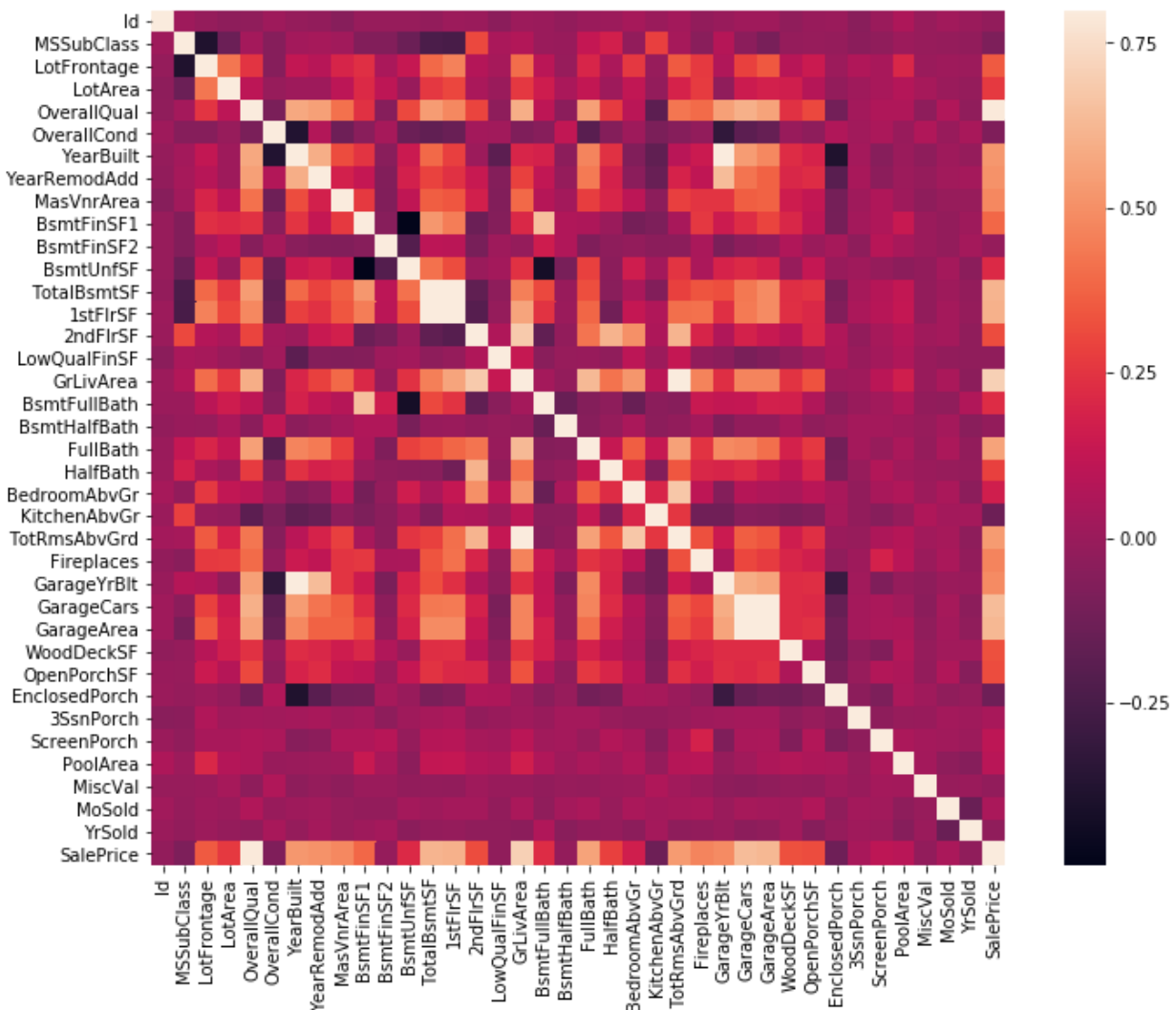


Skewness: 1.882876

Kurtosis: 6.536282

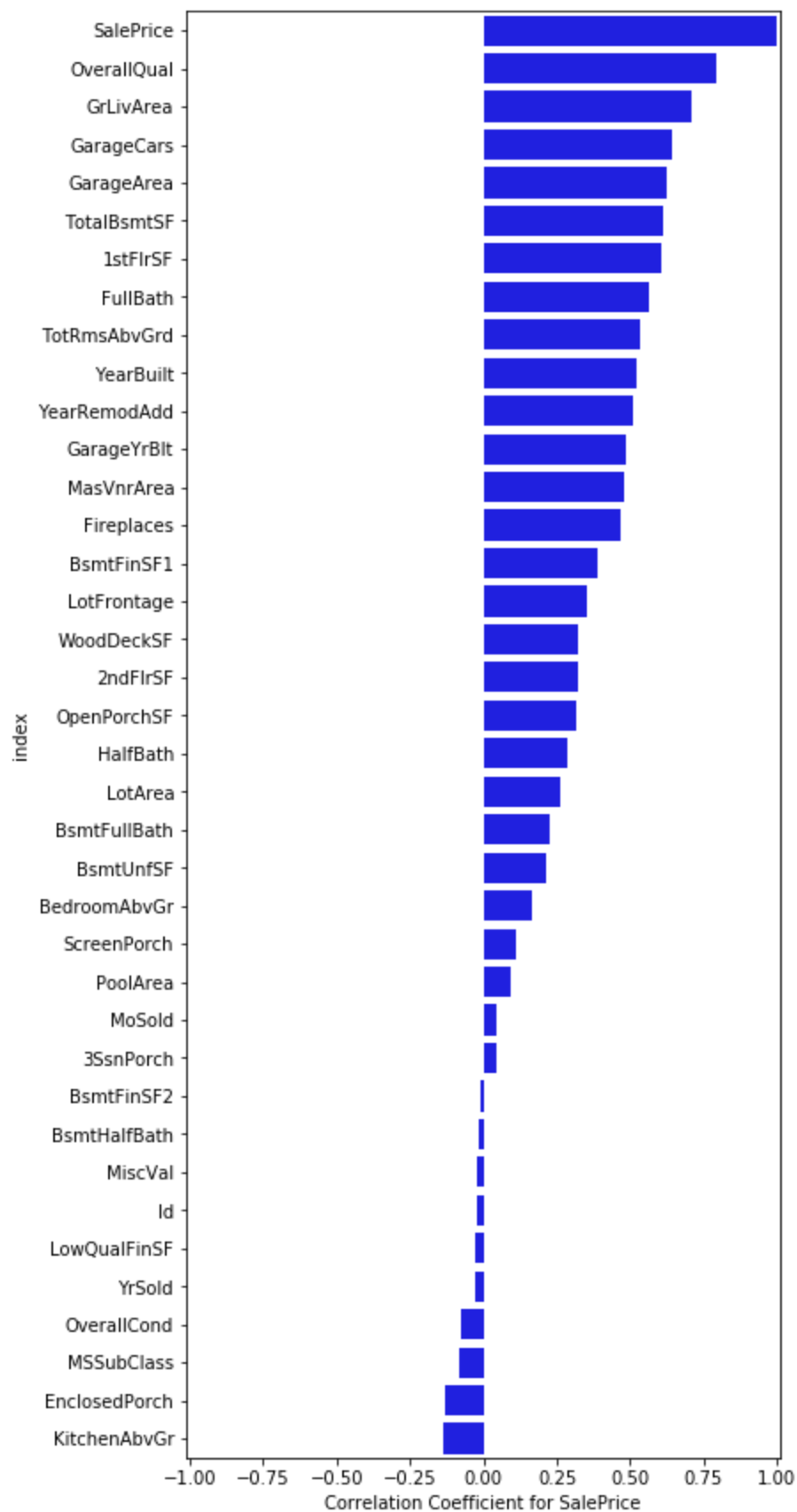
As shown by the plot above, the distribution of sale price is skewed right. Since the distribution is not totally normal, we chose to use a non-parametric approach, Kruskal-Wallis, as one of our methods.

Correlation analysis and the heatmaps for numeric and ordinal-categorical variables



The heatmap is the best way to display the relationship between every pair of variables. There are several off-diagonal white squares in this heatmap. For example, between the variables “GarageYrBlt” and “YearBuilt”. Another is between “GarageCars” and “GarageArea”. This is because there is a very high correlation within each of these pairs, as is suggested by intuition (e.g., often times, an old house has an old garage). Essentially, the information the first item in a given one of these pairs provides is almost the same as the information provided by the second item in the same given pair. So we can expect to find a high presence of multicollinearity at the sites of these pairs. Notice how the bottom edge and the right edge of the heatmap stand out by being lighter in hue compared to their surroundings. The squares in these locations represent correlations between our explanatory variables and our dependent variable. The higher these correlations are, the better. In this case, the heatmap shows an encouraging trend.

For the dependent variable “SalePrice”, we can see strong relationships between it and “OverallQual”, “GrLivArea”, “GarageCars”, “TotalBsmtSF”, and “YearBuilt”, among others.



The correlation coefficients for "SalePrice" are represented in the bargraph above. Clearly, "OverallQual", "GrLivArea", "GarageCars", "GarageArea", "TotalBsmtSF", "1stFirSF", "FullBath", "TotRmsAbvGrd", "YearBuilt", "YearRemodAdd", "GarageYrBit", "MasVnrArea", "MasVnrArea", "Fireplaces", "BsmtFinSF1", "LotFrontage", "WoodDeckSF", "2ndFlrSF", "OpenPorchSF", "HalfBath", "LotArea", "BsmtUnfSF", "BedroomAbvGr", "ScreenPorch", "PoolArea", "MoSold" and "3SsnPorch" have positive relationship. "BsmtHalfBath", "MiscVal", "Id", "LowQualFinSF", "YrSold", "OverallCond", "MSSubClass", "EnclosedPorch" and KitchenAbvGr" are negative correlation. (Id is negligible as an explanatory variable because it is only an arbitrary index introduced in the data generating phase.)

Kruskal-Wallis Hypothesis Test for non-ordinal categorical variable

Null hypothesis (same for each of the categorical variables we test here):

H_0 : All distributions are equal

Alternative hypothesis:

H_1 : At least one distribution is different from one other distribution

Since the sale price is not a normal distribution, a non parametric test is appropriate to test whether there is significant difference among groupings of SalePrice by the individual categories of each given categorical variable. There are three non-ordinal categorical variables where there is no significant difference in sale price. These variables are: Street, LandSlope, and Condition2.

II. Method:

1. Regression/Lasso:

After eliminating the categorical features that did not influence price using the Kruskal-Wallis H Test, we run some models to try to estimate price with the remaining features.

First, we utilized lasso (we checked that the features removed by lasso were somewhat correlated to the remaining ones, this is because lasso is inclined to reduce the coefficient to 0, so it is important to make sure that features that may not influence the model in a great proportion but are important are not removed). After running the model, we did cross validation to get the square root of the mse, which was 31258.7858\$.

The parameters we used for Lasso were:(incrementing to more iterations did not give a better result)

alpha = 2, max_iter = 5000, tol = 1e-20

Then, we utilize Recursive Feature Elimination (or RFECV), which works by recursively removing attributes and building a model on those attributes that remain. It uses the model accuracy to identify which attributes (and combinations of attributes) contribute the most to the predicting the target. We utilized the lasso coefficients from the previous model and ran the RFECV, which gave us our optimal number of features. Then, we tested the selected model with selected with our training data on our test data set and got 29170.01272 as our RMSE.

This shows us the coefficients for the optimal number of features after RFECV and Lasso:

```
File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

In [72]: print(RMSE(Df, namesfeatures)[0], RMSE(Df, namesfeatures)[1])
          ##using only the features recursive feature elimination after lasso we get a RMSE of:
          29170.01272105987 [(('MSSubClass', -122.53683692154652), ('Alley', 2836.890300515983), ('LotShape', -2506.648312378157
          7), ('OverallQual', 7324.418184421028), ('OverallCond', 5396.090875214641), ('YearBuilt', 214.49251762941697), ('Exte
          rQual', 8056.936687509088), ('ExterCond', -2610.088730326294), ('BsmQual', 5824.203649086903), ('BsmCond', -3275.31
          40198325677), ('BsmExposure', 6358.770994735396), ('BsmFinTypel', 1342.0840161197289), ('HeatingQC', 839.2423776474
          51), ('CentralAir', 581.2761803593459), ('1stFlrSF', 62.543545375216446), ('2ndFlrSF', 55.96433681083469), ('BsmFlt
          Bath', 8260.76257325948), ('BsmHalfBath', 5857.926110961587), ('FullBath', 2302.2325461211126), ('HalfBath', 3408.68
          4522277407), ('BedroomAbvGr', -5136.65092568845), ('KitchenAbvGr', -17616.644040729243), ('KitchenQual', 6695.753722
          276749), ('TotHwmAbvGr', 3622.886742462732), ('Functional', 6236.355346091791), ('Fireplacew', 9468.362879121565),
          ('FireplaceQu', -1894.7927019739852), ('GarageType', -2529.649939610165), ('GarageFinish', 1106.5850598755076), ('Gar
          ageCars', 9493.693112574634), ('GarageQual', 6493.810199546465), ('GarageCond', -8301.287461239412), ('PavedDrive', 9
          99.307061300188), ('3sanPorch', 35.879165115828584), ('PoolArea', 175.80381294505318), ('PoolQC', -17162.38856540653
          4), ('Fenceo', -376.21563934572106), ('MoSold', -845.9825775052095), ('MSZoning_C (all)', -25897.533582000367), ('MSZo
          ning_FV', 367.7909191696325), ('MSZoning_RM', 5044.98216317707), ('MSZoning_RM', -165.20178842453436), ('LandContour
          _Bnk', -9708.610100823878), ('LandContour_HLS', 5584.883960737843), ('LandContour_Low', -4637.007269253987), ('LotCon
          fig_Corner', 970.8364701019277), ('LotConfig_CuldSac', 10472.160009886536), ('LotConfig_FR2', -1907.989323122394),
          ('LotConfig_FR3', -22468.605771503615), ('Neighborhood_Blmngtn', 1708.5085934241952), ('Neighborhood_Blueste', 0.0),
          ('Neighborhood_BrDale', 27451.896045537498), ('Neighborhood_BrkSide', 3315.0869790584316), ('Neighborhood_ClearCr', 6
          798.242084009963), ('Neighborhood_CollCr', -7684.727409609204), ('Neighborhood_Crawfor', 17748.450047742725), ('Neig
          hborhood_Edwards', -12800.151495386019), ('Neighborhood_Gilbert', -13258.614217711644), ('Neighborhood_IDOTRR', -507
          2.689965457473), ('Neighborhood_MeadowV', 9569.69905418303), ('Neighborhood_Mitchel', -5997.724000393883), ('Neighbo
          hood_NAmes', -6673.86229647055), ('Neighborhood_NPKVill', 32693.283795267358), ('Neighborhood_NWAmes', -8426.39914464
          3978), ('Neighborhood_Nokridge', 25208.984428211188), ('Neighborhood_Nridgnt', 44462.767632529765), ('Neighborhood_Old
          Town', -8730.352351878952), ('Neighborhood_SMIU', -3852.382128051), ('Neighborhood_SawyerW', -3222.1086130309413),
          ('Neighborhood_Somerst', 5344.223501707743), ('Neighborhood_StoneBr', 57238.007691426195), ('Neighborhood_Timber', -2
          293.8162599012644), ('Neighborhood_Veenker', 8871.05405309793), ('Conditionl_Feedr', -567.2666480746271), ('Condition
          l_Norm', 1182.654006852006), ('Conditionl_PosN', -19164.759973679147), ('Conditionl_RRAe', -1185.572462706003), ('C
          conditionl_RRAn', 12255.98105285728), ('Conditionl_RRNe', -8303.907390800658), ('BldgType_lFam', 3337.0911836793784),
          ('BldgType_2fmCon', 15295.662490433937), ('BldgType_Twnhs', -24083.056188793846), ('BldgType_TwnhSh', -18225.92838245
          908), ('RoofStyle_Flat', -0.0), ('RoofStyle_Gable', -1185.4626539341868), ('RoofStyle_Gambrel', 0.0), ('RoofStyle_Hi
          p', 5454.805177483883), ('RoofStyle_Mansard', 6847.025808316998), ('RoofStyle_Shed', -4719.035801867924), ('RoofMatl_
          ClyTile', -438075.156084648), ('RoofMatl_Membran', 32341.28496828473), ('RoofMatl_Metal', 0.0), ('RoofMatl_Roll', -11
          695.970291523265), ('RoofMatl_TarGrv', -12770.555307576367), ('RoofMatl_WdShake', -19970.402080933454), ('RoofMatl_W
          dShngl', 99425.90004574393), ('Exterior1st_AsbShng', 9930.665673491692), ('Exterior1st_AspShm', 0.0), ('Exterior1st_
          BrkFace', 14679.54687429952), ('Exterior1st_Cemmntd', 28056.40185604104), ('Exterior1st_HdBoard', -5410.01188357477
          4), ('Exterior1st_InsStucc', -6621.890363462178), ('Exterior1st_MetalSd', 11344.133369886835), ('Exterior1st_Plywood',
          -5039.35668167419), ('Exterior1st_Stucco', 10288.409437591321), ('Exterior1st_VinylSd', 2533.6590820923875), ('Exteri
          or1st_WdShing', 3040.6183317092837), ('Exterior1st_AsbShng', -4476.570648852368), ('Exterior1st_AspShm', 250.7833006
          7944524), ('Exterior1st_Brk cmn', -6842.0005858930409), ('Exterior1st_Cemmntd', -22027.688761235902), ('Exterior1st_Hd
          Board', 1109.912161302637), ('Exterior1st_InsStucc', 5881.008599503434), ('Exterior1st_MetalSd', -4072.5004280615744),
          ('Exterior1st_Other', -31609.51281309596), ('Exterior1st_Plywood', 0.0), ('Exterior1st_Stone', -23350.44719632034),
          ('Exterior1st_Stucco', -10424.314888004423), ('Exterior1st_VinylSd', -601.3479571718484), ('Exterior1st_Wd Shng', 208
          8.488352195024), ('Exterior1st_Wd Shng', -4692.8590700272625), ('MasVnrType_BrCmn', -4085.583197833488), ('MasVnrType
          eNone', -246.6323409661847), ('MasVnrType_Stone', 4466.902710968929), ('Foundation_BrkTil', -6125.250563014455), ('F
          oundation_PConc', 277.67421842929), ('Foundation_Slab', 8736.829070636404), ('Foundation_Wood', -23077.83343826827
          3), ('Heating_GasA', 1295.676574863671), ('Heating_GasW', 6431.221873714051), ('Heating_OthW', -40574.61507199289),
          ('Heating_Wall', 18423.68919103293), ('Electrical_FuseA', -4544.805725915716), ('Electrical_FuseF', -5821.04260889932
          5), ('Electrical_SBrkr', -7983.80249661505), ('MiscFeature_Gar2', -6629.149723021818), ('MiscFeature_Othr', 23022.071
          699002627), ('MiscFeature_Shed', 779.3623233869221), ('MiscFeature_TenC', -75240.47731961719), ('SaleType_COD', -1374
          5.407628140132), ('SaleType_Con', 17385.87507934808), ('SaleType_ConLw', -13362.6501248543), ('SaleType_New', 42452.
          599997984646), ('SaleType_WD', -8956.101665585069), ('SaleCondition_AdjLand', 7851.671222857444), ('SaleCondition_All
          oca', 9840.24754154221), ('SaleCondition_Family', -2504.4106906482434), ('SaleCondition_Normal', 7298.485645638147),
          ('SaleCondition_Partial', -27875.022650948224), ('LotFrontage', 38.4505293631606), ('Housestyle', -7215.26532067416
          4)]

In [73]: def RFECVcvvvvvvvvilinearreg(Df, Features_used):
```

Another approach we utilized with RFECV was linear regression. Our RMSE was 29824.8428.

When trying to decide which of these two models to utilize, we took into consideration that the method with lasso and RFECV has some penalty if more coefficients (variables) are in the model. This penalty helps us avoid overfitting the model and it still works well with a different data set.

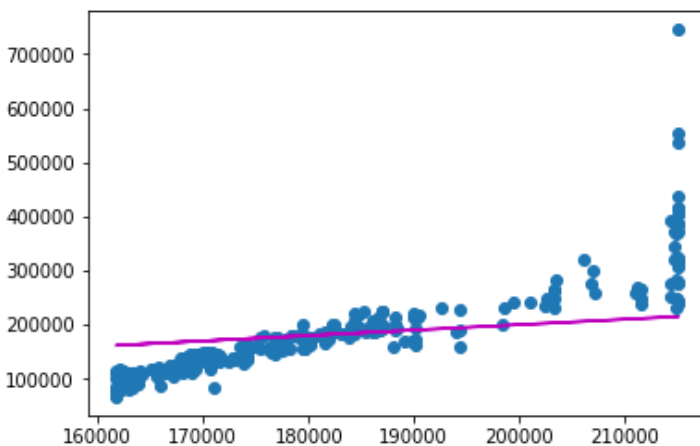
2. LightGBM

DataSet preparation: We divided the dataset into three parts. Training and testing datasets for model building. A true testing dataset is used for comparing algorithms.

Parameter setting: There are four major parameters in LightGBM: learning_rate, sub_feature, num_leaves, min_data, and max_depth. Our initial setting and the corresponding MSE were:

```
params['learning_rate'] = 0.003
params['sub_feature'] = 0.5
params['num_leaves'] = 10
params['min_data'] = 50
params['max_depth'] = 10
MSE (before optimization)= 4764732533.46014$
```

The prediction plot shows the overall prediction performance in the true testing dataset:



y-axis: True SalePrice ; x-axis:

Prediction

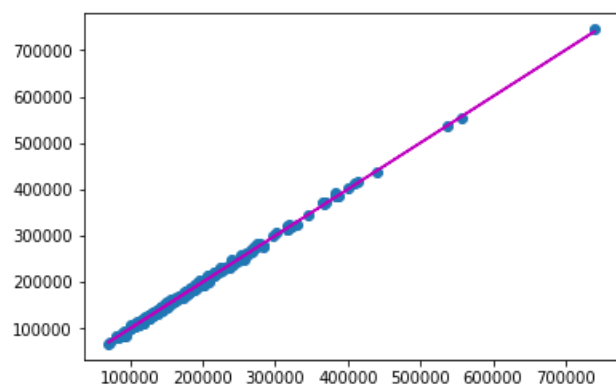
After using the above parameters to optimize the coefficients on the train data, we found the values to replace the above parameters with that minimized the MSE on the test data.

in during the modeling process using training and testing datasets.

```
params['learning_rate'] = 0.121
params['sub_feature'] = 0.3
params['num_leaves'] = 8
params['min_data'] = 2
params['max_depth'] = 5
```

Then we used these parameters to predict the sale price in the true testing dataset. Our MSE for that LightGBM is 9819810.0.

The prediction plot shows the overall prediction performance in the true testing



dataset:

References :

<https://www.kaggle.com/yanpapadakis/houseprices-eda> anova

<https://www.kaggle.com/sxliuliang9494/xgboost-lasso> xgboost

<https://www.kaggle.com/pmarcelino/comprehensive-data-exploration-with-python> Correlation