

```

1 import numpy as np
2 import pandas as pd
3 from matplotlib import pyplot as plt
4 import seaborn as sns
5 from sklearn.preprocessing import LabelEncoder, StandardScaler
6 from sklearn.model_selection import train_test_split
7 from sklearn.linear_model import LinearRegression
8 from sklearn import metrics
9 from sklearn import ensemble
10 from sklearn.linear_model import Lasso,Ridge
11 #load train data
12 df_data=pd.read_csv("/content/data_price.csv")
13 df_data.head()

```

```

↳
   Id  MSSubClass  MSZoning  LotFrontage  LotArea  Street  Alley  LotShape  LandContour  Utilities  LotConfig  LandSlop
0   1           60        RL           65.0    8450   Pave   NaN        Reg           Lvl        AllPub      Inside      C
1   2           20        RL           80.0    9600   Pave   NaN        Reg           Lvl        AllPub       FR2      C
2   3           60        RL           68.0   11250   Pave   NaN        IR1           Lvl        AllPub      Inside      C
3   4           70        RL           60.0    9550   Pave   NaN        IR1           Lvl        AllPub    Corner      C
4   5           60        RL           84.0   14260   Pave   NaN        IR1           Lvl        AllPub       FR2      C

```

5 rows x 81 columns

```

1 #to know each and every column execute the following
2 print(df_data.columns)
3 print(df_data.shape)

```

```

↳ Index(['Id', 'MSSubClass', 'MSZoning', 'LotFrontage', 'LotArea', 'Street',
        'Alley', 'LotShape', 'LandContour', 'Utilities', 'LotConfig',
        'LandSlope', 'Neighborhood', 'Condition1', 'Condition2', 'BldgType',
        'HouseStyle', 'OverallQual', 'OverallCond', 'YearBuilt', 'YearRemodAdd',
        'RoofStyle', 'RoofMatl', 'Exterior1st', 'Exterior2nd', 'MasVnrType',
        'MasVnrArea', 'ExterQual', 'ExterCond', 'Foundation', 'BsmtQual',
        'BsmtCond', 'BsmtExposure', 'BsmtFinType1', 'BsmtFinSF1',
        'BsmtFinType2', 'BsmtFinSF2', 'BsmtUnfSF', 'TotalBsmtSF', 'Heating',
        'HeatingQC', 'CentralAir', 'Electrical', '1stFlrSF', '2ndFlrSF',
        'LowQualFinSF', 'GrLivArea', 'BsmtFullBath', 'BsmtHalfBath', 'FullBath',
        'HalfBath', 'BedroomAbvGr', 'KitchenAbvGr', 'KitchenQual',
        'TotRmsAbvGrd', 'Functional', 'Fireplaces', 'FireplaceQu', 'GarageType',
        'GarageYrBlt', 'GarageFinish', 'GarageCars', 'GarageArea', 'GarageQual',
        'GarageCond', 'PavedDrive', 'WoodDeckSF', 'OpenPorchSF',
        'EnclosedPorch', '3SsnPorch', 'ScreenPorch', 'PoolArea', 'PoolQC',
        'Fence', 'MiscFeature', 'MiscVal', 'MoSold', 'YrSold', 'SaleType',
        'SaleCondition', 'SalePrice'],
        dtype='object')
(1460, 81)

```

```

1 total = df_data.isnull().sum().sort_values(ascending=False)
2 percent = (df_data.isnull().sum()/df_data.isnull().count()).sort_valu
3 missing_data = pd.concat([total, percent], axis=1, keys=['Total', 'Pe
4 missing_data.head(20)

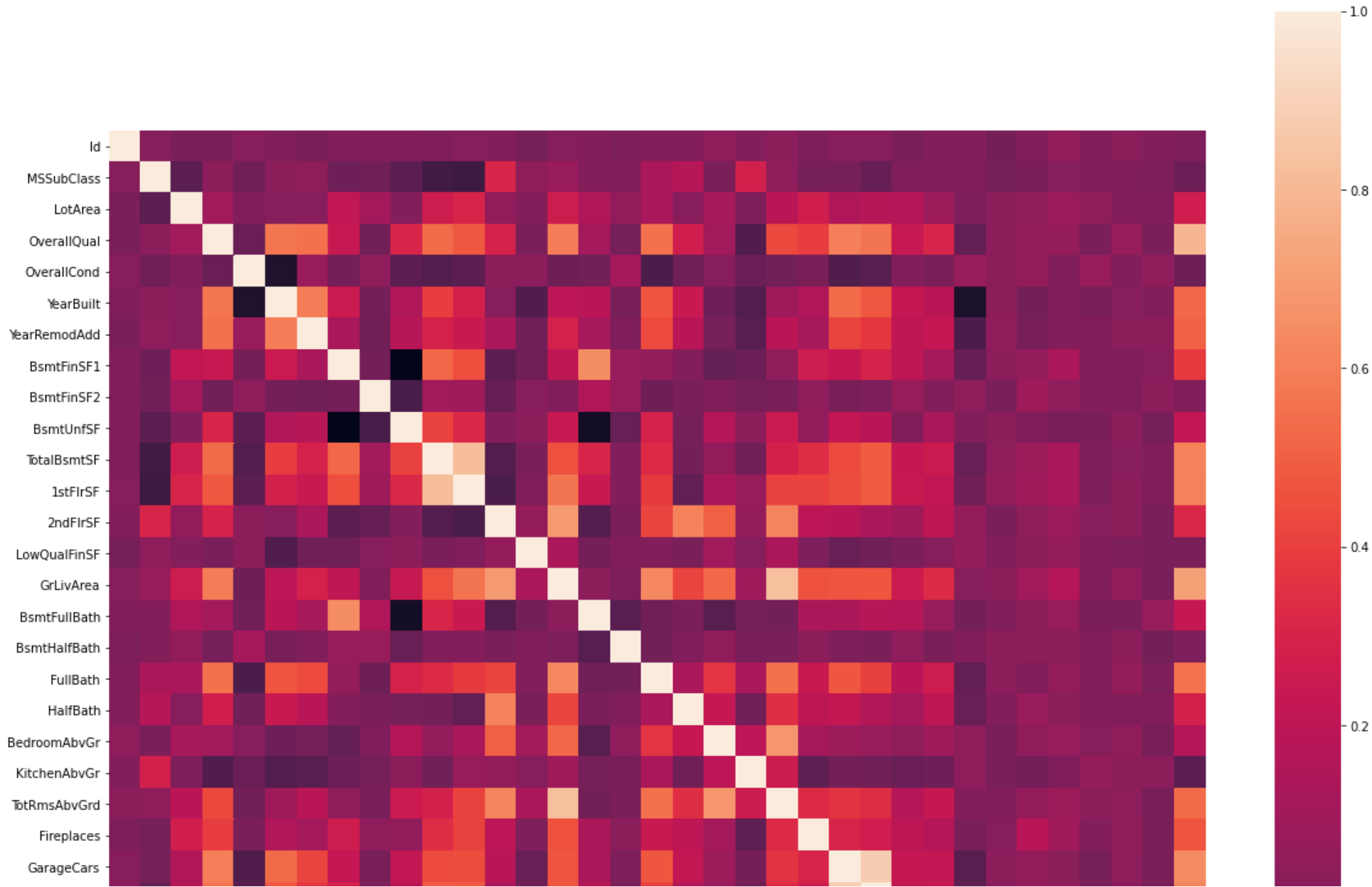
```

↳

	Total	Percent
PoolQC	1453	0.995205
MiscFeature	1406	0.963014
Alley	1369	0.937671
Fence	1179	0.807534
FireplaceQu	690	0.472603
LotFrontage	259	0.177397
GarageCond	81	0.055479

```
1 df_data= df_data.drop(missing_data[missing_data['Total']>1].index.val
2 df_data= df_data.drop(df_data.loc[df_data['Electrical'].isnull()].ind
    GarageFinish      81  0.055479
1 corr_mat=df_data.corr()
2 fi,ax=plt.subplots(figsize=(20,20))
3 sns.heatmap(corr_mat,square=True)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f7428024630>
```



```
1 del df_data['Id']

1 le=LabelEncoder()
2 cat_mask= df_data.dtypes=='object'
3 cat_cols= df_data.columns[cat_mask].tolist()
4 cat_cols
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f7428024630>
```



```

['MSZoning',
 'Street',
 'LotShape',
 'LandContour',
 'Utilities',
 'LotConfig',
 'LandSlope',
 'Neighborhood',
 'Condition1',
 'Condition2',
 'BldgType',
 'HouseStyle',

1  #Lets convert the columns to one hot encoding
2  df_data[cat_cols]=df_data[cat_cols].apply(lambda x: le.fit_transform(x
3  df_data_c = df_data.copy()
4  #get_dummies is used for one hot encoding
5  df_data_c = pd.get_dummies(df_data_c,columns=cat_cols)

    'Heating'

1  x_train, x_test, y_train, y_test = train_test_split(df_data_c.drop('S
2  y_train= y_train.values.reshape(-1,1)
3  y_test= y_test.values.reshape(-1,1)

    'Functional',

```

Normalize the values in train and test using Standard Scaler function.

```

    'SaleCondition' ]

1  sc_X = StandardScaler()
2  sc_y = StandardScaler()
3  x_train = sc_X.fit_transform(x_train)
4  x_test = sc_X.fit_transform(x_test)
5  y_train = sc_X.fit_transform(y_train)
6  y_test = sc_y.fit_transform(y_test)

1  lm = LinearRegression()
2  lm.fit(x_train,y_train)
3  #predictions on train data
4  x_pred = lm.predict(x_train)
5  x_pred = x_pred.reshape(-1,1)
6  #Prediction of validation data
7  y_predictions = lm.predict(x_test)
8  y_predictions= y_predictions.reshape(-1,1)
9  def scores_(y,x):
10     print('MAE:', metrics.mean_absolute_error(y, x))
11     print('MSE:', metrics.mean_squared_error(y, x))
12     print('RMSE:', np.sqrt(metrics.mean_squared_error(y, x)))
13     print('R2 Score:' ,metrics.r2_score(y,x))
14  print('InSample_accuracy')
15  scores_(y_train, x_pred)
16  print('-----')
17  print('OutSample_accuracy')
18  scores_(y_test,y_predictions)

[> InSample_accuracy
MAE: 0.1828459849312525
MSE: 0.08080438090379947
RMSE: 0.2842611139494804
R2 Score: 0.9191956190962005
-----
OutSample_accuracy
MAE: 29345417710.666924
MSE: 3.034299649072542e+21
RMSE: 55084477387.66832
R2 Score: -3.034299649072542e+21

```

The model performed really well on training data with a good 0.92 r2 score and < 1 RMSE score but with test data the performance is nowhere near good. This is a clear overfitting model. Reason might be because of numerous features. To tackle this we can perform Ridge and Lasso regularization.

```

1  gularization(model,alpha_range):
2  se_score_insample=[]
3  se_score_outsample=[]
4  _score_insample=[]
5  _score_outsample=[]
6  for i in alpha_range:
7     regularization = model(alpha=i,normalize=True)
8     regularization.fit(x_train,y_train)
9     y_pred_train = regularization.predict(x_train)
10    y_pred_train = y_pred_train.reshape(-1,1)

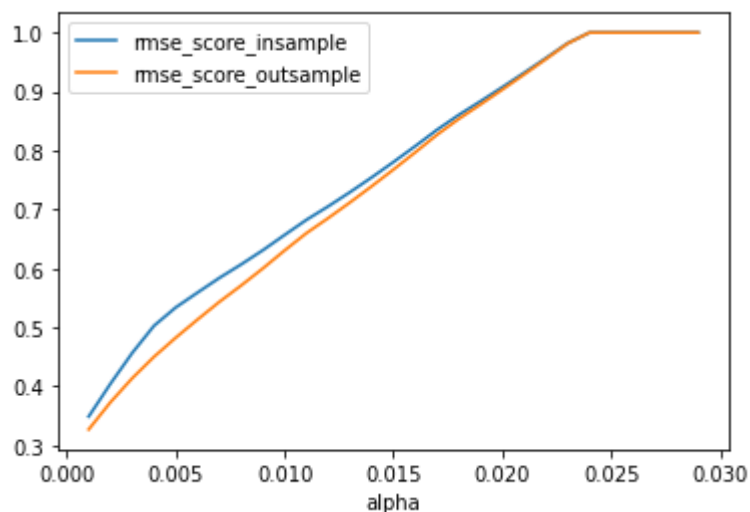
```

```

11 y_pred_test=regularization.predict(x_test)
12 y_pred_test = y_pred_test.reshape(-1,1)
13 rmse_score_insample.append(np.sqrt(metrics.mean_squared_error(y_tra
14 rmse_score_outsample.append(np.sqrt(metrics.mean_squared_error(y_te
15 r2_score_insample.append(metrics.r2_score(y_train, y_pred_train))
16 r2_score_outsample.append(metrics.r2_score(y_test, y_pred_test))
17 =pd.DataFrame()
18 ['alpha']=alpha_range
19 ['rmse_score_insample'] = rmse_score_insample
20 ['rmse_score_outsample']= rmse_score_outsample
21 ['r2_score_insample'] = r2_score_insample
22 ['r2_score_outsample'] = r2_score_outsample
23 turn df.plot(x = 'alpha', y = ['rmse_score_insample', 'rmse_score_out
24 range_lasso = np.arange(0.001,0.03,0.001)
25 regularization(Lasso,alpha_range_lasso))

```

➤ AxesSubplot(0.125,0.125;0.775x0.755)



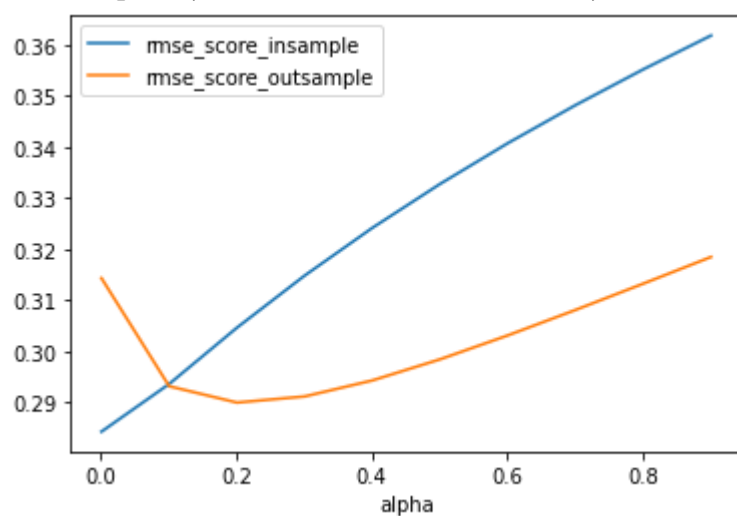
We can see that there is no huge difference in in sample and out sample RMSE scores so Lasso has resolved overfitting. One observation here is that after  $\alpha = 0.017$  there is no difference in RMSE scores of In sample and Out sample.

```

1 alpha_range_ridge = np.arange(0.001,1,0.1)
2 print(regularization(Ridge,alpha_range_ridge))
3 #writing functions helps reduce redundant lines of code as seen #abov

```

➤ AxesSubplot(0.125,0.125;0.775x0.755)



We see in the graph that around  $\alpha = 0.1$  there is no much difference in the RMSE scores and clearly there is no sign of over fitting as there is very less difference of insample and outsample RMSE scores as compared to huge difference in Linear Regression.