## Question 1

What is the optimal value of alpha for ridge and lasso regression? What will be the changes in the model if you choose double the value of alpha for both ridge and lasso? What will be the most important predictor variables after the change is implemented?

The optimal value of alpha for ridge and lasso regression

- Ridge Alpha 1
- lasso Alpha 10

**Ridge Regression**

```
#Change the alpha value from 1 to 2

alpha = 3

ridge2 = Ridge(alpha=alpha)

ridge2.fit(X_train1, y_train)

Ridge(alpha=3)

# Lets calculate some metrics such as R2 score, RSS and RMSE

y_pred_train = ridge2.predict(X_train1)

y_pred_test = ridge2.predict(X_test1)


metric2 = []

r2_train_lr = r2_score(y_train, y_pred_train)

print(r2_train_lr)

metric2.append(r2_train_lr)


r2_test_lr = r2_score(y_test, y_pred_test)

print(r2_test_lr)

metric2.append(r2_test_lr)


rss1_lr = np.sum(np.square(y_train - y_pred_train))

print(rss1_lr)
```

```
metric2.append(rss1_lr)


rss2_lr = np.sum(np.square(y_test - y_pred_test))

print(rss2_lr)

metric2.append(rss2_lr)


mse_train_lr = mean_squared_error(y_train, y_pred_train)

print(mse_train_lr)

metric2.append(mse_train_lr**0.5)


mse_test_lr = mean_squared_error(y_test, y_pred_test)

print(mse_test_lr)

metric2.append(mse_test_lr**0.5)


#Alpha 1

#R2score(train) 0.884340040460635

#R2score(test)  0.869613280468847

0.8797315810932456

0.8710282148272899

607995142958.1411

320928407278.46216

680845624.8131479

729382743.8146868
```

R2score on training data has decreased but it has increased on testing data


## Lasso

```
#Changed alpha 10 to 20

alpha =20

lasso20 = Lasso(alpha=alpha)

lasso20.fit(X_train1, y_train)

Lasso(alpha=20)
```

```python
# Lets calculate some metrics such as R2 score, RSS and RMSE
y_pred_train = lasso20.predict(X_train1)
y_pred_test = lasso20.predict(X_test1)


metric3 = []
r2_train_lr = r2_score(y_train, y_pred_train)
print(r2_train_lr)
metric3.append(r2_train_lr)


r2_test_lr = r2_score(y_test, y_pred_test)
print(r2_test_lr)
metric3.append(r2_test_lr)


rss1_lr = np.sum(np.square(y_train - y_pred_train))
print(rss1_lr)
metric3.append(rss1_lr)


rss2_lr = np.sum(np.square(y_test - y_pred_test))
print(rss2_lr)
metric3.append(rss2_lr)


mse_train_lr = mean_squared_error(y_train, y_pred_train)
print(mse_train_lr)
metric3.append(mse_train_lr**0.5)


mse_test_lr = mean_squared_error(y_test, y_pred_test)
print(mse_test_lr)
metric3.append(mse_test_lr**0.5)


#R2score at alpha-10
#0.8859222400899005
```

#0.8646666084570094

0.8854019697956436

0.8670105921065014

579329522996.7144

330925704432.26794

648745266.5136778

752103873.7096999

R2score of training data has decrease and it has increase on testing data


#important predictor variables

betas = pd.DataFrame(index=X_train1.columns)

betas.rows = X_train1.columns

betas['Ridge2'] = ridge2.coef_

betas['Ridge'] = ridge.coef_

betas['Lasso'] = lasso.coef_

betas['Lasso20'] = lasso20.coef_

pd.set_option('display.max_rows', None)

betas.head(68)

- LotArea---------------Lot size in square feet
- OverallQual---------Rates the overall material and finish of the house
- OverallCond--------Rates the overall condition of the house
- YearBuilt-------------Original construction date
- BsmtFinSF1--------Type 1 finished square feet
- TotalBsmtSF------- Total square feet of basement area
- GrLivArea-----------Above grade (ground) living area square feet
- TotRmsAbvGrd----Total rooms above grade (does not include bathrooms)
- Street_Pave--------Pave road access to property
- RoofMatl_Metal----Roof material_Metal

Predictors are same but the coefficient of this predictor has changed

Question 2

You have determined the optimal value of lambda for ridge and lasso regression during the assignment. Now, which one will you choose to apply and why?

The r2_score of lasso is slightly higher than lasso for the test dataset so we will choose lasso regression to solve this problem.


Question 3

After building the model, you realised that the five most important predictor variables in the lasso model are not available in the incoming data. You will now have to create another model excluding the five most important predictor variables. Which are the five most important predictor variables now?

X_train

y_train

X_train1.columns

Index(['LotArea', 'OverallQual', 'OverallCond', 'YearBuilt', 'BsmtFinSF1', 'TotalBsmtSF', '1stFlrSF', '2ndFlrSF', 'GrLivArea', 'BedroomAbvGr', 'TotRmsAbvGrd', 'Street_Pave', 'LandSlope_Sev', 'Condition2_PosN', 'RoofStyle_Shed', 'RoofMatl_Metal', 'Exterior1st_Stone', 'Exterior2nd_CBlock', 'ExterQual_Gd', 'ExterQual_TA', 'BsmtCond_Po', 'KitchenQual_TA', 'Functional_Maj2', 'SaleType_CWD', 'SaleType_Con'], dtype='object')

LotArea,OverallQual,YearBuilt,BsmtFinSF1,TotalBsmtSF are the top 5 important predictors.

Let's drop these columns

X_train2 = X_train1.drop(['LotArea','OverallQual','YearBuilt','BsmtFinSF1','TotalBsmtSF'],axis=1)

X_test2 = X_test1.drop(['LotArea','OverallQual','YearBuilt','BsmtFinSF1','TotalBsmtSF'],axis=1)

X_train2.head()


## Lasso

```
# alpha 10
alpha =10
lasso21 = Lasso(alpha=alpha)
lasso21.fit(X_train2, y_train)
Lasso(alpha=10)
# Lets calculate some metrics such as R2 score, RSS and RMSE
y_pred_train = lasso21.predict(X_train2)
y_pred_test = lasso21.predict(X_test2)

metric3 = []
r2_train_lr = r2_score(y_train, y_pred_train)
print(r2_train_lr)
metric3.append(r2_train_lr)

r2_test_lr = r2_score(y_test, y_pred_test)
```

```
print(r2_test_lr)
metric3.append(r2_test_lr)

rss1_lr = np.sum(np.square(y_train - y_pred_train))
print(rss1_lr)
metric3.append(rss1_lr)

rss2_lr = np.sum(np.square(y_test - y_pred_test))
print(rss2_lr)
metric3.append(rss2_lr)

mse_train_lr = mean_squared_error(y_train, y_pred_train)
print(mse_train_lr)
metric3.append(mse_train_lr**0.5)

mse_test_lr = mean_squared_error(y_test, y_pred_test)
print(mse_test_lr)
metric3.append(mse_test_lr**0.5)
#R2score at alpha-10
#0.8859222400899005
#0.8646666084570094
0.7988346707068132
0.758810320925813
1016954777102.8657
600167078819.8159
1138807141.2126155
1364016088.2268543
R2score of training and testing data has decreased

#important predictor variables
betas = pd.DataFrame(index=X_train2.columns)
betas.rows = X_train1.columns
betas['Lasso21'] = lasso21.coef_
pd.set_option('display.max_rows', None)
betas.head(68)
```

five most important predictor variables

- 11stFlrSF-----------First Floor square feet
- GrLivArea-----------Above grade (ground) living area square feet
- Street_Pave---------Pave road access to property
- RoofMatl_Metal------Roof material_Metal
- RoofStyle_Shed------Type of roof(Shed)

Question 4

How can you make sure that a model is robust and generalisable? What are the implications of the same for the accuracy of the model and why?

The model should be generalized so that the test accuracy is not lesser than the training score. The model should be accurate for datasets other than the ones which were used during training. Too much importance should not given to the outliers so that the accuracy predicted by the model is high. To ensure that this is not the case, the outliers analysis needs to be done and only those which are relevant to the dataset need to be retained. Those outliers which it does not

make sense to keep must be removed from the dataset. If the model is not robust, It cannot be trusted for predictive analysis.