

## 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?

**Answer:**

The optimal values for Ridge and Lasso regression using top 120 variables selected by RFE are :

- Ridge Alpha : 5.0
- Lasso Alpha : 0.001

As we increase the value of alpha for ridge and lasso the model will become simpler i.e., the bias will increase and variance will decrease

As we increase the value of  $\lambda$  then the magnitude of the coefficients decreases.

Most Important Predictor Variables Before Change :

First Half Variables are the sorted using absolute ridge coefficients

Second Half Variables are the sorted using absolute lasso coefficients

	Ridge	Lasso
Neighborhood_NoRidge	0.539098	0.592742
OverallQual_9	0.494384	0.830237
OverallQual_10	0.468908	0.833928
FullBath_3	0.465469	0.530678
BsmtQual_TA	0.367518	0.373752
TotRmsAbvGrd_11	0.347359	0.505044
Neighborhood_NridgHt	0.344950	0.352086
BsmtQual_Fa	0.331857	0.371310
Neighborhood_Somerst	0.313458	0.393006
Fireplaces_3	0.306188	0.442772

```
Index(['Neighborhood_NoRidge', 'OverallQual_9', 'OverallQual_10', 'FullBath_3',  
      'BsmtQual_TA', 'TotRmsAbvGrd_11', 'Neighborhood_NridgHt', 'BsmtQual_Fa',  
      'Neighborhood_Somerst', 'Fireplaces_3'],  
      dtype='object')
```

	Ridge	Lasso
OverallQual_10	0.468908	0.833928
OverallQual_9	0.494384	0.830237
Neighborhood_NoRidge	0.539098	0.592742
FullBath_3	0.465469	0.530678
TotRmsAbvGrd_11	0.347359	0.505044
Fireplaces_3	0.306188	0.442772
OverallQual_8	0.182730	0.411655
Neighborhood_Somerst	0.313458	0.393006
BsmtQual_TA	0.367518	0.373752
BsmtQual_Fa	0.331857	0.371310

```
Index(['OverallQual_10', 'OverallQual_9', 'Neighborhood_NoRidge', 'FullBath_3',  
      'TotRmsAbvGrd_11', 'Fireplaces_3', 'OverallQual_8',  
      'Neighborhood_Somerst', 'BsmtQual_TA', 'BsmtQual_Fa'],  
      dtype='object')
```

Most Important Predictor Variables After Change :

First Half Variables are the sorted using absolute ridge coefficients

Second Half Variables are the sorted using absolute lasso coefficients

```

                Ridge      Lasso
Neighborhood_NoRidge  0.446595  0.506984
OverallQual_9         0.430324  0.840617
FullBath_3            0.393552  0.496926
OverallQual_10        0.378727  0.816993
BsmtQual_TA           0.337015  0.331769
Neighborhood_NridgHt  0.314769  0.301451
BsmtExposure_Gd       0.287828  0.276494
1stFlrSF              0.285616  0.273562
KitchenQual_TA        0.282519  0.250262
BsmtQual_Fa           0.260388  0.291325
Index(['Neighborhood_NoRidge', 'OverallQual_9', 'FullBath_3', 'OverallQual_10',
      'BsmtQual_TA', 'Neighborhood_NridgHt', 'BsmtExposure_Gd', '1stFlrSF',
      'KitchenQual_TA', 'BsmtQual_Fa'],
      dtype='object')

                Ridge      Lasso
OverallQual_9         0.430324  0.840617
OverallQual_10        0.378727  0.816993
Neighborhood_NoRidge  0.446595  0.506984
FullBath_3            0.393552  0.496926
OverallQual_8          0.177479  0.420404
TotRmsAbvGrd_11       0.252978  0.355101
BsmtQual_TA           0.337015  0.331769
GarageCars_3          0.255577  0.322583
Neighborhood_Somerst  0.245755  0.308089
Neighborhood_NridgHt  0.314769  0.301451
Index(['OverallQual_9', 'OverallQual_10', 'Neighborhood_NoRidge', 'FullBath_3',
      'OverallQual_8', 'TotRmsAbvGrd_11', 'BsmtQual_TA', 'GarageCars_3',
      'Neighborhood_Somerst', 'Neighborhood_NridgHt'],
      dtype='object')

```

Metric	R2 Score (Train)	R2 Score (Test)	RSS (Train)	RSS (Test)	MSE (Train)	MSE (Test)
Linear Regression	0.9	0.84	102.25	72.88	0.32	0.41
Ridge Regression	0.89	0.85	110.0	67.41	0.33	0.39
Lasso Regression	0.89	0.85	109.62	67.73	0.33	0.39
Ridge Regression Double Lambda	0.89	0.85	117.35	67.29	0.34	0.39
Lasso Regression Double Lambda	0.88	0.85	119.47	67.5	0.34	0.39

## 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?

**Answer:**

Metric	R2 Score (Train)	R2 Score (Test)	RSS (Train)	RSS (Test)	MSE (Train)	MSE (Test)
Linear Regression	0.9	0.84	102.25	72.88	0.32	0.41
Ridge Regression	0.89	0.85	110.0	67.41	0.33	0.39
Lasso Regression	0.89	0.85	109.62	67.73	0.33	0.39

The optimal values for Ridge and Lasso regression using top 120 variables selected by RFE are :

- Ridge Alpha : 5.0
- Lasso Alpha : 0.001

Difference in r2 values for train and test data have improved from **0.6 to 0.4**.

The RSS for the test data has reduced from **72.88 to 67.41 and 67.73** for Ridge and Lasso respectively (lower the value better)

The MSE for the test data has reduced from **72.88 to 67.41 and 67.73** for Ridge and Lasso respectively (lower the value better)

Also, the Lasso performs feature elimination by making their coefficient 0 indirectly making the model simple, robust and generalized in nature.

```
In [271]: #sum of coefficients
betas.sum()
```

```
Out[271]: Linear    5.481599
Ridge      2.490063
Lasso      3.873862
dtype: float64
```

```
In [275]: #Number of variables in model after feature elimination
betas[betas!=0].count()
```

```
Out[275]: Linear    121
Ridge      118
Lasso       81
dtype: int64
```

In the given data set Lasso eliminated 3 features from 120 selected by RFE.

In the given data set Lasso eliminated 40 features from 120 selected by RFE.

Clearly the Lasso Regularization will produce a simple model and yield better results on unseen data.

### 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?

**Answer:**

<b>Metric</b>	<b>R2 Score (Train)</b>	<b>R2 Score (Test)</b>	<b>RSS (Train)</b>	<b>RSS (Test)</b>	<b>MSE (Train)</b>	<b>MSE (Test)</b>
<b>Linear Regression</b>	0.9	0.84	102.25	72.88	0.32	0.41
<b>Ridge Regression</b>	0.89	0.85	110.0	67.41	0.33	0.39
<b>Lasso Regression</b>	0.89	0.85	109.62	67.73	0.33	0.39
<b>Ridge Regression Double Lambda</b>	0.89	0.85	117.35	67.29	0.34	0.39
<b>Lasso Regression Double Lambda</b>	0.88	0.85	119.47	67.5	0.34	0.39
<b>Ridge Regression Drop Top 5</b>	0.87	0.84	135.18	74.23	0.36	0.41
<b>Lasso Regression Drop Top 5</b>	0.87	0.84	132.13	72.51	0.36	0.41

There is a drop in r2 values for train and test data

There is increase in RSS and MSE values for train and test data

The Top 5 predictors after changes

Using Ridge :

- 'OverallQual\_5'
- 'OverallQual\_6',
- 'OverallQual\_4'
- 'KitchenQual\_TA',
- 'KitchenQual\_Fa'

Using Lasso :

- 'OverallQual\_5'
- 'OverallQual\_4'
- 'OverallQual\_6'
- 'OverallQual\_3'
- 'Fireplaces\_3'

#### 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?

**Answer:**

We can use L1 and L2 Regularization to make sure that the model is robust and generalisable.

A model is considered to be robust if its output is consistently accurate even if one or more of the input variables changed.

A model is considered to be generalisable if it is able to predict on the complete population using a model trained by a subset of the data

Simpler models are usually more 'generic'.

Simpler models are more robust.

Typical ways of looking the complexity of a model.

1. Number of parameters required to specify the model completely.
2. The degree of the function, if it is a polynomial
3. Size of the best-possible representation of the model.
4. The depth or size of a decision tree.

The difference in  $r^2$  for train and test is small

The variance and bias should be as low as possible (Bias Variance Trade-off)

The model should not be complex i.e., the degree of polynomial or depth of the tree should be as small

**Implications:**

- The accuracy of the model on train set might go down
- The accuracy of the model on test set might go up

