

**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 to double the value of alpha for both ridge and lasso? What will be the most important predictor variables after the change is implemented?

**Answer:**

**Ridge regression:** Optimal value of alpha obtained for ridge regression is 10.0.

**Lasso regression:** Optimal value of alpha obtained for lasso regression is .001

As the value of alpha increases, there is a significant bias and variance trade-off is observed. As the value of alpha (lambda) increases it will try to shrink the co-efficient towards zero which will make y predicted value less i.e., " $Y = B_0 + B_1 * X$ " and hence RSS will increase.

As the change is not significant though  $r^2\_score$  of the model is dropped due to an increase in RSS.

**The most significant parameter in lasso regression is "GrLivArea", where as**

**The most significant parameter in ridge regression is "OverallQual".**

**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:**

Lasso and ridge regression both are used for regularization. Lasso regression helps to perform feature selection in model by shrinking co-efficient towards zero as well as makes some of the coefficient zero. Ridge regression is preferred when all variables have same level of significance, and it shrinks the coefficients towards zero not zero.

In our analysis lasso and ridge both have similar results, **but lasso results tend to seem better** as  $r^2\_score$  and RSS, RMSE results both are same for training and test data whereas for ridge regression we can notice slight variation in results.

In addition, from EDA we found that not all variables have same level of impact on response variable and only few has higher co-efficient, **we choose lasso over ridge** as it will help the model in feature selection.

**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:**

Most important features obtained for the lasso regression before were:

GrLivArea
Sale Type (type: New)
OverallQual
Neighborhood_Crawfor
Overall Cond

After dropping the five most important features:

Neighbourhood (Type: NridgHt)
Neighbourhood (Type: StoneBr)
Central Air (Type: Y)
2ndFlrSF
1stFlrSF

**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:**

Model tends to perform better if it has low variance and low bias, and it does not overfit/underfit the data. Model over-fitting can be observed by analysing the  $r^2_{score}$ , comparing RSS & RMSE in training and test data. For ex: In housing price prediction we clearly saw that  $r^2_{score}$  for training data was approx. 90 and for test data it was  $< 0$ . It was a sign of over fitting of the model.

In addition, model residuals should be plotted against Y and there should not be any variance observed in the plot. If a variance is observed in histogram of error or residual plots against Y, it may indicate non-linearity of the data and a model may be overfitted. To solve the case polynomial regression, data-transformation and non-linear regression can be used.

To make a sure a model is robust and generalisable, model should have less features in final model. In addition,  $r^2_{score}$  for training and test data should be same. K fold cross validation should be used to find out the hidden patten in underlying data and a best model should be chosen. Model is considered robust and generalisable if it has less variance and tends to give similar result on test data