## 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?**

Optimal Value of lambda for Ridge : 100
Optimal Value of lambda for Lasso : 0.001

After doubling the values of alpha for both Ridge and Lasso :
- R2 score using Ridge (Train) remained almost same (~0.94)
- R2 score using Ridge (Test) increased from 0.9015 to 0.9419
- R2 score using Lasso (Train) decreased from 0.9473 to .8998
- R2 score using Lasso (Test) increased from 0.8937 to 0.9118

- The most significant variables in Ridge are GrLivArea, TotalBsmtSF, 1stFlrSF, OverallQual_8, BsmtFinSF1
- The most significant variables in Lasso are GrLivArea, TotalBsmtSF, OverallQual_8, OverallQual_9, BsmtFinSF1

## 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?**

Both the options gave good scores. Lasso seems to be a good option to go with since it refines the parameter selection and the coefficients aren't too large to choose from.

## Question 3

**After building the model, you realized 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?**

After dropping the first 5 top predictors and re building the model, here are the top 5 predictors:
**Lasso** : 1stFlrSF, 2ndFlrSF, BsmtFinSF1, BsmtUnfSF, LotArea
**Ridge** : 1stFlrSF, 2ndFlrSF, BsmtFinSF1, GarageArea, LotArea

**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?**

A model is robust and generalisable when any variation or new data or unseen data does not affect the performance of the model. These can be achieved by making sure that the model does not overfit. In an overfit model, all the noise data is taken into account and this is very high variance and very small change in data will lead to drastic change in the model.
On the same lines, a model that is very complex will have high accuracy. To make the model more robust and generalizable, we will have to decrease variance which will lead to some bias. Addition of bias means that accuracy will decrease. This has to be taken care by striking some balance between model accuracy and complexity. This can be achieved by Regularization techniques like Ridge Regression and Lasso.