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?

Ans: the optimal value of alpha ridge regression value is 2 and Lasso regression value is 0.0001

```python
In [55]: # Building Ridge Model by doubling the value of alpha to 4
         ridge_double = Ridge(alpha=4,random_state=100)
         ridge_double.fit(X_train_rfe2,y_train)
         ridge_double_coef = ridge_double.coef_
         y_test_pred = ridge_double.predict(X_test_rfe2)
         print('The R2 Score of the model on the test dataset for doubled alpha is',r2_score(y_test, y_test_pred))
         print('The MSE of the model on the test dataset for doubled alpha is', mean_squared_error(y_test, y_test_pred))
         ridge_double_coeff = pd.DataFrame(np.atleast_2d(ridge_double_coef),columns=X_train_rfe2.columns)
         ridge_double_coeff = ridge_double_coeff.T
         ridge_double_coeff.rename(columns={0: 'Ridge Doubled Alpha Co-Efficient'},inplace=True)
         ridge_double_coeff.sort_values(by=['Ridge Doubled Alpha Co-Efficient'], ascending=False,inplace=True)
         print('The most important predictor variables are as follows:')
         ridge_double_coeff.head(20)

         The R2 Score of the model on the test dataset for doubled alpha is 0.8259998671982055
         The MSE of the model on the test dataset for doubled alpha is 0.001862290533613281
         The most important predictor variables are as follows:
```

Out[55]:

| | Ridge Doubled Alpha Co-Efficient |
|---|---|
| Total_sqr_footage | 0.149028 |
| GarageArea | 0.091803 |

When we double the value of alpha for our ridge regression no we will take the value of alpha equal to 10 the model will apply more penalty on the curve and try to make the model more generalized that is making model more simpler and no thinking to fit every data of the data set .from the graph we can see that when alpha is 10 we get more error for both test and train. Similarly when we increase the value of alpha for lasso we try to penalize more our model and more coefficient of the variable will reduced to zero, when we increase the value of our r2 square also decreases.

## 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?

Ans: It is important to regularize coefficients and improve the prediction accuracy also with the decrease in variance, and making the model interpretably.

Ridge regression, uses a tuning parameter called lambda as the penalty is square of magnitude of coefficients which is identified by cross validation. Residual sum or squares can be small by using the penalty. The penalty is lambda times sum of squares of the coefficients, hence the coefficients that have greater values gets penalized. As we increase the value of lambda the variance in model is dropped and bias remains constant.

Lasso regression, uses a tuning parameter called lambda as the penalty is absolute value of magnitude of coefficients which is identified by cross validation. As the lambda value increases Lasso shrinks the coefficient towards zero and it make the variables exactly equal to 0. Lasso also does variable selection. When lambda value is small it performs simple linear regression and as lambda value increases, shrinkage takes place and variables with 0 value are neglected by the model

**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?

Ans: these 5 most important predictor variables that will be excluded are :-

1. Total_sqr_footage 2. GarageArea 3. TotRmsAbvGrd

4. OverallCond 5. LotArea

```
#Removing the 5 most important predictor variables from the incoming dataset
X_test_rfe3 = X_test_rfe2.drop(['Total_sqr_footage','GarageArea','TotRmsAbvGrd','OverallCond','LotArea'],axis=1)
X_train_rfe3 = X_train_rfe2.drop(['Total_sqr_footage','GarageArea','TotRmsAbvGrd','OverallCond','LotArea'],axis=1)

# Building Lasso Model with the new dataset
lasso3 = Lasso(alpha=0.0001,random_state=100)
lasso3.fit(X_train_rfe3,y_train)
lasso3_coef = lasso3.coef_
y_test_pred = lasso3.predict(X_test_rfe3)
print('The R2 Score of the model on the test dataset is',r2_score(y_test, y_test_pred))
print('The MSE of the model on the test dataset is', mean_squared_error(y_test, y_test_pred))
lasso3_coeff = pd.DataFrame(np.atleast_2d(lasso3_coef),columns=X_train_rfe3.columns)
lasso3_coeff = lasso3_coeff.T
lasso3_coeff.rename(columns={0: 'Lasso Co-Efficient'},inplace=True)
lasso3_coeff.sort_values(by=['Lasso Co-Efficient'], ascending=False,inplace=True)
print('The most important predictor variables are as follows:')
lasso3_coeff.head(5)
```

```
The R2 Score of the model on the test dataset is 0.7330077964268464
The MSE of the model on the test dataset is 0.0028575670906482538
The most important predictor variables are as follows:
```

| | Lasso Co-Efficient |
|---|---|
| LotFrontage | 0.146535 |
| Total_porch_sf | 0.072445 |
| House Style_2.5Unf | 0.062900 |
| House Style_2.5Fin | 0.050487 |
| Neighborhood_Veenker | 0.042532 |

## 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?.

Ans: 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.