## Assignment 2

### Question 1:

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:

➤ Alpha is the regularization parameter that controls the amount of regularization applied in the models. Higher values of alpha result in stronger regularization, which can help prevent over fitting but may also lead to under fitting.

 

The optimal values of alpha for ridge are as below:
➤
➤
➤

```
In [276]:  # Printing the best hyperparameter alpha
           print(model_cv.best_params_)

{'alpha': 100}
```

The optimal values of alpha for Lasso are as below:
➤
➤
➤

```
In [284]:  # Printing the best hyperparameter alpha
           print(model_cv.best_params_)

{'alpha': 1000}
```
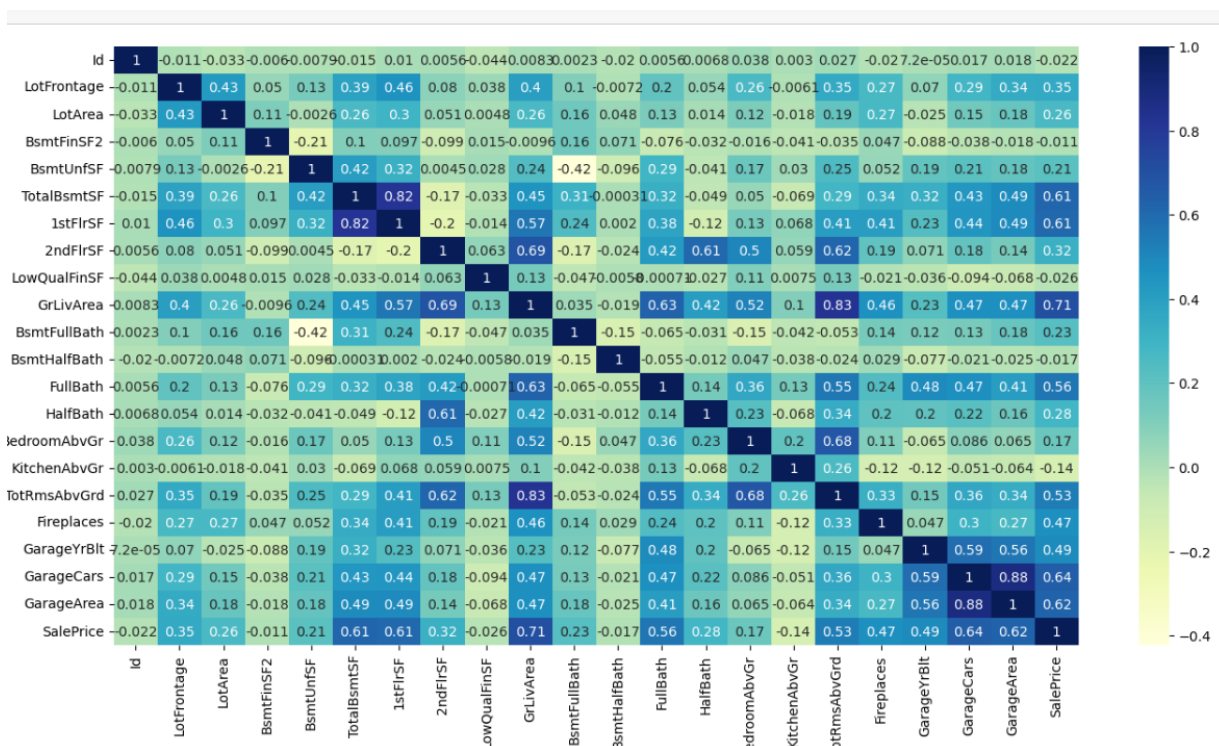
➤ If we double the value of alpha for both ridge and lasso regression, the models will be more heavily regularized. This means that the models will put stronger constraints on the coefficients and may shrink them further towards zero. Consequently, the impact on the model will be:

Ridge Regression: Increasing alpha in ridge regression will increase the amount of shrinkage applied to the coefficients. The coefficients will be further reduced towards zero but will not be completely eliminated. This can result in a simpler model with potentially fewer important predictors, as some coefficients may be further diminished or become close to zero.

Lasso Regression: Doubling alpha in lasso regression will lead to increased sparsity in the model. The higher regularization will push more coefficients to exactly zero, effectively removing them from the model. As a result, the number of important predictor variables will decrease, and the model will become more simplified with a reduced number of active predictors.

➤ The most important predictor variable is in data set 'Saleprice'.

## 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:**

lambda for ridge and lasso regression during the assignment is as below:

Out[299]:

| | Metric | Linear Regression | Ridge Regression | Lasso Regression |
|---|---|---|---|---|
| 0 | R2 Score (Train) | 7.086886e-01 | 7.086881e-01 | 7.086630e-01 |
| 1 | R2 Score (Test) | 7.359387e-01 | 7.358763e-01 | 7.358549e-01 |
| 2 | RSS (Train) | 1.858770e+12 | 1.858773e+12 | 1.858933e+12 |
| 3 | RSS (Test) | 7.443129e+11 | 7.444888e+11 | 7.445492e+11 |
| 4 | MSE (Train) | 4.266777e+04 | 4.266781e+04 | 4.266965e+04 |
| 5 | MSE (Test) | 4.122311e+04 | 4.122798e+04 | 4.122965e+04 |

**Analysis:**
R2 Score:
➤ The R2 score measures the proportion of the variance in the target variable that is predictable from the predictors.
➤ For both ridge and lasso regression, the R2 scores for both the train and test sets are very close to the R2 score of linear regression. This suggests that ridge and lasso regression have achieved a similar level of predictive performance.
➤ Based on the R2 score alone, there is no significant advantage of using ridge or lasso regression over linear regression.

Residual Sum of Squares (RSS):
➤ The RSS measures the total squared difference between the predicted and actual target values.

➢ Both ridge and lasso regression have slightly higher RSS values compared to linear regression for both the train and test sets.

➢ Linear regression has the lowest RSS values, indicating a better fit to the training and test data compared to ridge and lasso regression.

Mean Squared Error (MSE):

> ➢ The MSE measures the average squared difference between the predicted and actual target values.
>
> ➢ Similar to the RSS, both ridge and lasso regression have slightly higher MSE values compared to linear regression for both the train and test sets.
>
> ➢ Linear regression has the lowest MSE values, indicating a better fit to the training and test data compared to ridge and lasso regression.

Based on the analysis, linear regression seems to perform slightly better than ridge and lasso regression in terms of the RSS and MSE on both the train and test sets. the trade-off between simplicity and predictive performance. Ridge and lasso regression offer benefits in scenarios where multicollinearity or variable selection is important.

Therefore, As our primary goal is predictive performance, we can choose linear regression. However, As per our concerns about multicollinearity or the interpretability of the model, we can consider ridge regression as it offers regularization to mitigate multicollinearity while maintaining most of the predictors. Lasso regression is also suitable as we suspect that only a subset of predictors is truly important and wish to perform variable 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:**

The most important predictor variables are as below:
 'TotalBsmtSF','1stFlrSF', 'GarageArea', 'TolRmsAbvGrd' and 'GarageCars'.

**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:**

To ensure that a model is robust and generalizable, you can take the following steps:

1. Train/Test Split: we should Split dataset into separate training and testing subsets. The training set is used to train the model, while the testing set is used to evaluate its

performance on unseen data. This helps assess how well the model generalizes to new, unseen observations.

2. Cross-Validation: We have to Utilize techniques like k-fold cross-validation to further assess the model's performance. Cross-validation involves splitting the data into multiple subsets or folds, training the model on different combinations of these folds, and evaluating its performance. This helps estimate the model's performance on unseen data more reliably and reduces the dependency on a single train/test split.

3. Evaluate Multiple Metrics: Use a range of evaluation metrics to assess the model's performance, such as accuracy, precision, recall, F1-score, or mean squared error. Considering multiple metrics provides a comprehensive view of the model's strengths and weaknesses, enabling better understanding of its generalizability.

4. Regularization: We have to apply regularization techniques like ridge regression or lasso regression to control overfitting. Regularization helps prevent the model from fitting the noise in the training data too closely, promoting better generalization to unseen data.

5. Feature Engineering and Selection: We have to engage in proper feature engineering and feature selection techniques to ensure that the model focuses on the most relevant and informative predictors. Removing irrelevant or redundant features can enhance the model's generalizability by reducing noise and overfitting.

6. Test on Diverse Data: We have to assess the model's performance on diverse subsets of the data, ensuring that it performs consistently well across different scenarios. This includes testing the model on data collected at different times, from different sources, or representing different subgroups of the population.

The implications of ensuring model robustness and generalizability for model accuracy are as follows:

1. Lower Over fitting: By focusing on generalization, the model aims to reduce overfitting, which occurs when the model fits the training data too closely and performs poorly on new data. This may result in a decrease in the training set accuracy. However, the model is expected to have better accuracy on unseen data, which is more important for real-world applications.

2. Improved Performance on Unseen Data: By following robust modeling practices, the model's accuracy on unseen data, such as the testing set or new data samples, is expected to be more reliable. This enables the model to make more accurate predictions in real-world scenarios.

3. Reduced Bias: Ensuring that the model is robust and generalizable helps reduce bias in the model's predictions. A bias occurs when the model consistently underestimates or overestimates the target variable. A robust model with good generalization properties is more likely to provide unbiased predictions.

4. Enhanced Model Stability: A robust and generalizable model is less sensitive to small variations in the data. This means that even if the training data changes slightly, the model's performance is expected to be consistent and stable, resulting in improved reliability.

Overall, focusing on model robustness and generalizability may lead to a small reduction in accuracy on the training set but is essential for ensuring better performance and reliability on unseen data, which is crucial for real-world applications and decision-making.