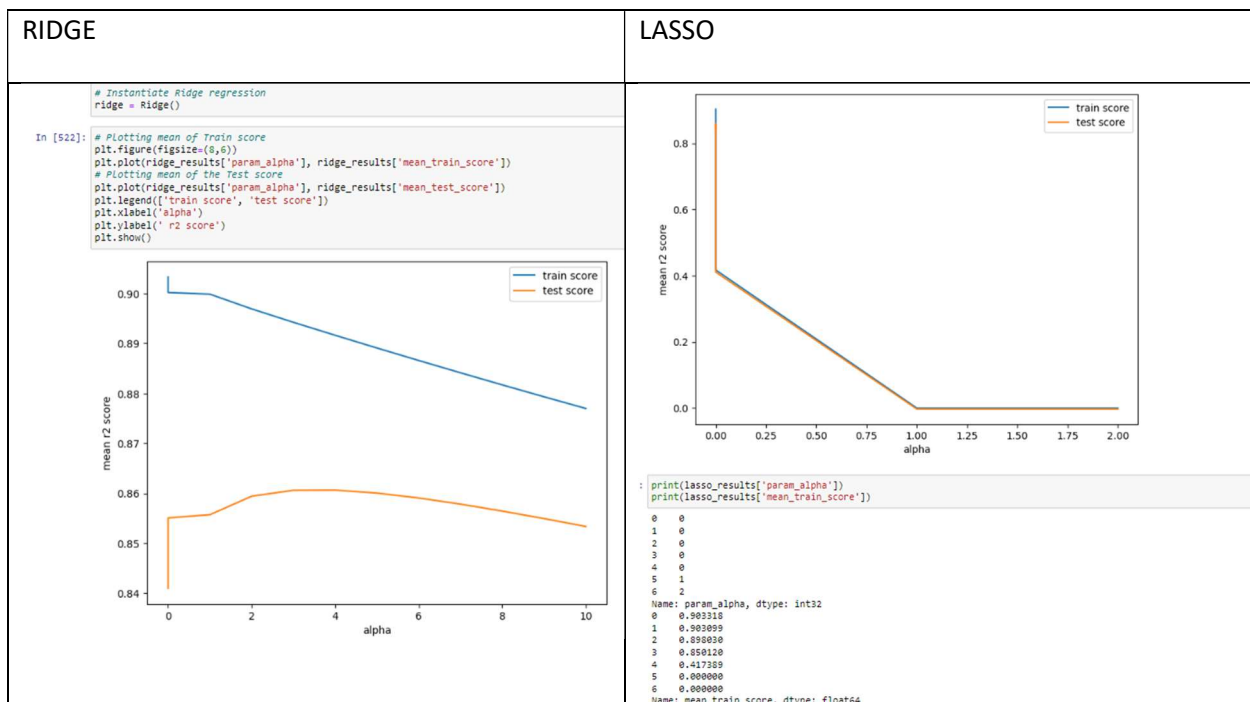# Assignment Part-II

**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

Basis the assignment of house price prediction, optimal value of alpha for ridge and lasso regression is below.

| RIDGE | LASSO |
|---|---|
|  |  |

the graph says lower the value of lambda (close to 0) , higher the accuracy. here it is coincidentally same

*Train Score*

r2 score goes down when the lambda value increase means the error is increasing. Model is overfitting and generalised.At 0.002 (close to 0) the train set accury is highest(more than 80%).

*Test Score*

At alpha = 0.002 the test accuracy is highest (more than 80%).

After alpha=0.002, the r2 score started decreasing as the alpha is increasing. Select the value of aplha for which the test score peaks up. In our case at alpha=0.002, the error is less in the test set and so accuracy is more close to 80%.

So, the optimum alpha will be 0.002, this level error and generalisation of the model for makes a simpler mode

Generally, smaller values of alpha result in models that are closer to ordinary linear regression, while larger values of alpha introduce more regularization and can help prevent overfitting.

**Double the value of alpha for both Ridge and Lasso, the following changes could occur:**

Ridge Regression:

The coefficients would be further shrunk towards zero, resulting in a more regularized model.

The model's complexity would decrease, and it would be less prone to overfitting.
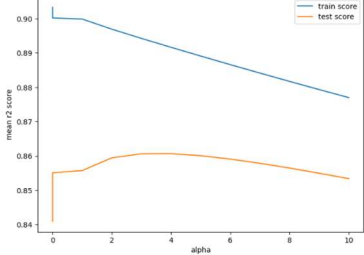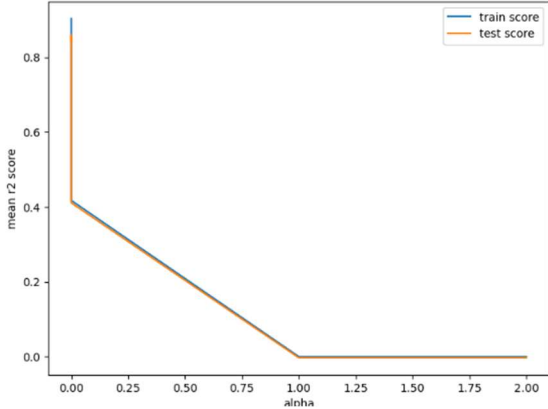
Lasso Regression:

More coefficients would be driven to exactly zero, leading to even sparser models.

| ### Top 10 features of Ridge regression and coefficients | ### Lasso Top 10 features and Co-Efficients |
|---|---|
| {'OverallQual': 0.1495, 'GrLivArea': 0.1407, 'OverallCond': 0.1008, 'GarageArea': 0.0838, 'FullBath': 0.0646, 'BsmtFullBath': 0.0477, 'Fireplaces': 0.0462, 'BedroomAbvGr': 0.0444, 'd_BsmtQual': 0.041, 'Neighborhood_StoneBr': 0.0373} | {'OverallQual': 0.2628, 'GarageArea': 0.0822, 'GrLivArea': 0.077, 'FullBath': 0.0565, 'Fireplaces': 0.0535, 'd_GarageFinish': 0.0423, 'd_HeatingQC': 0.0313, 'MSZoning_RL': 0.0295, 'd_BsmtFinType1': 0.0292, 'd_BsmtQual': 0.0235}¶ |

## 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?

| RIDGE | LASSO |
|---|---|
|  |  |
| Ridge Regression Model:<br><br> Train data:<br>R2 score: 0.8937809012072492<br>RMSE score: 0.03958866748353199<br><br> Test data:<br>R2 score: 0.8740921714800556<br>RMSE score: 0.04417766528950032<br><br><br>Conclusion We have a good train score 85% and good test score as well 83%. That means what the model learnt in the train set it performed well in the test set. | Lasso Regression Model<br><br> Train data:<br>R2 score: 0.8069457722085366<br>RMSE score: 0.05337147880821713<br><br> Test data :<br>R2 score: 0.800881671942512<br>RMSE score: 0.05555613339120312 |

- The R2 test score on the Ridge Regression Model is better than Lasso Regression Model.
- RMSE on the Ridge Regression Model is better than Lasso Regression Model.

## 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?

five most important predictor variables in the current lasso model

1. 'OverallQual': 0.2628,

2. 'GarageArea': 0.0822,

3. 'GrLivArea': 0.077,

4. 'FullBath': 0.0565,

5. 'Fireplaces': 0.0535,

Dropped five most important predictor variables  and performed lasso regression

```
R2 Score -test dataset 0.8320680460799884
 MSE - test dataset is 0.0026030716848734364
```
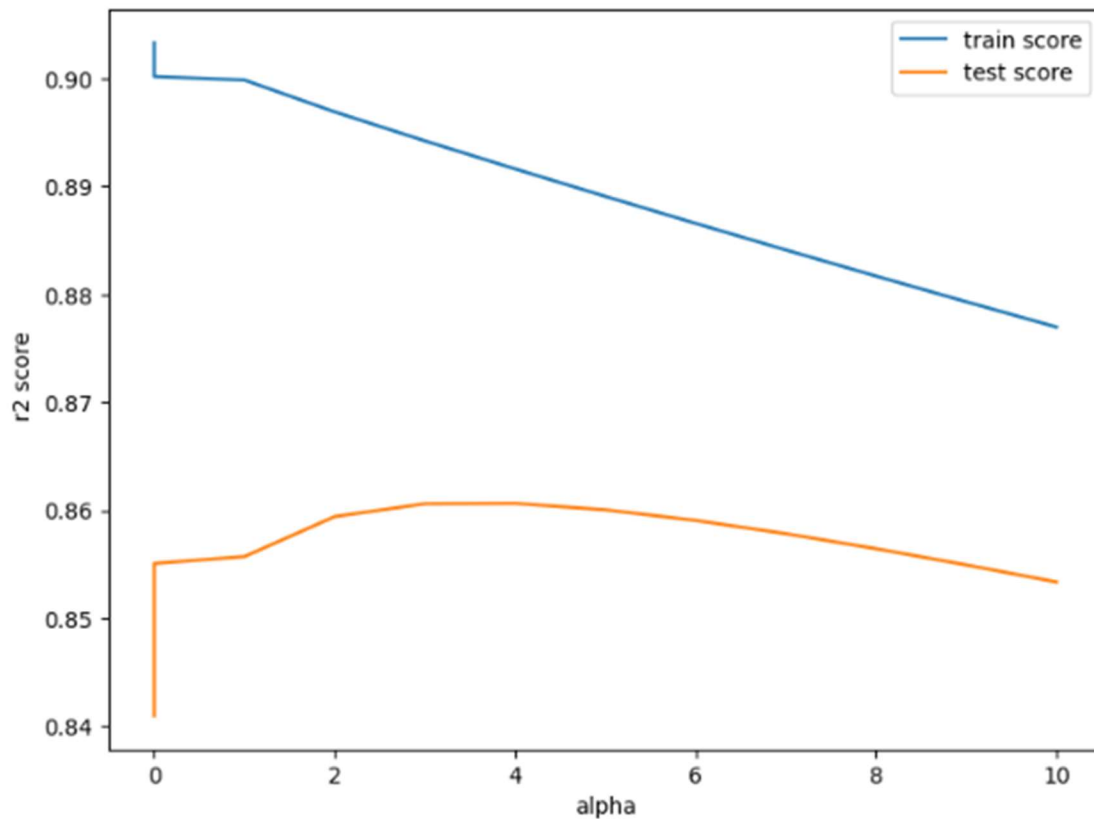
# The next set of 5  predictor variables are:

|  | lassoR_coef |
| --- | --- |
| TotalBsmtSF | 0.337070 |
| 2ndFlrSF | 0.168409 |
| OverallCond | 0.136727 |
| LotArea | 0.074614 |
| Neighborhood_StoneBr | 0.064120 |

**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?

1. **Regularization**: Apply regularization techniques like Lasso and Ridge to prevent overfitting. Regularization helps the model capture underlying patterns without fitting to noise in the training data.
2. **Hyperparameter Tuning**: Tune hyperparameters using techniques like grid search or random search. Avoid tuning based solely on the test set performance to prevent overfitting to the test set.



In our model, the r2 score started decreasing as the alpha is increasing after the value 2 and model accuracy started decreasing.

So, we pick the value of alpha =2, the error is much less in the test set and accuracy is good nearly 81%.

So, the optimum value for alpha could be 2, that will have the correct balance between the error and the generalization, to create a simple model.