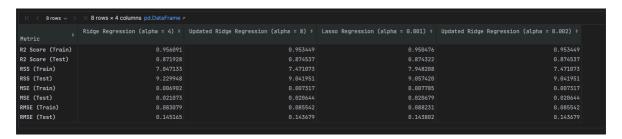
# **Housing Prices Case Study Assignment - Part 2**

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 -

Optimal value of lambda <u>for Ridge Regression</u> = 4 Optimal value of lambda <u>for Lasso</u> = 0.001

If we double the values of alpha for both Ridge and Lasso i.e. to 8 and 0.002 respectively, we didn't saw any major change in the R2 scores for Train and Test data set -



# 1. Ridge LinearRegression

- R2 Score for Train set remained pretty much same i.e. 0.95
- R2 Score for Test Set also remained same i.e. 0.87

## 2. Lasso LinearRegression

- R2 Score for Train set remained pretty much same i.e. 0.95
- R2 Score for Test Set also remained same i.e. 0.87

The most important predictor variables after the change is implemented will be

- 1. RoofMatl\_CompShg
- 2. GrLivArea
- 3. RoofMatl\_Tar&Grv
- 4. OverallQual
- 5. RoofMatl\_WdShngl
- 6. RoofMatl WdShake
- 7. OverallCond
- 8. GarageCars
- 9. RoofMatl\_Membran
- 10. RoofMatl Roll

and lasso regression during the assignment. Now, which one will you choose to apply and why?

### Answer -

The model we will choose to apply will depend on the use case.

- If we have too many variables and one of our <u>primary goal is feature</u> <u>selection</u>, then we will use **Lasso Linear Regression**.
- If we don't want to get too large coefficients and <u>reduction of</u> <u>coefficient magnitude is one of our primary goals</u>, then we will use **Ridge Linear Regression**.

Question 2 - 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-

We will have to drop the top 5 features in Lasso model and build the model again

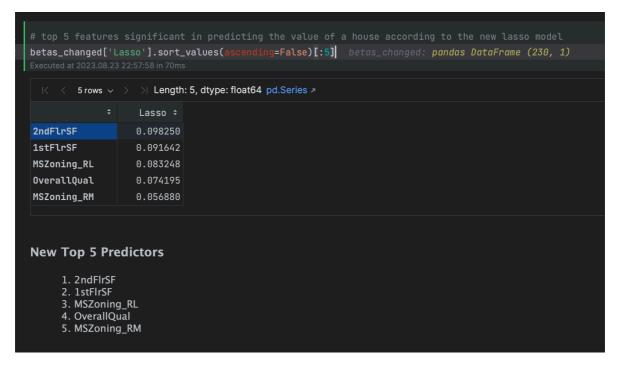
Top 5 features based on original optimal alpha values was -

Sl.No.	Feature	Description
1	RoofMatl_CompShg	Roof Material is "Standard (Composite) Shingle"
2	RoofMatl_Tar&Grv	Roof Material is "Gravel & Tar"
3	GrLivArea	Above grade (ground) living area square feet
4	RoofMatl_WdShngl	Roof Material is "Wood Shingles"
5	RoofMatl_WdShake	Roof Material is "Wood Shakes"

After we remove the above 5 columns and retrain the model using Lasso Regression, we get the following metrics -

```
R-Squared (Train) = 0.93
R-Squared (Test) = 0.87
RSS (Train) = 10.81
RSS (Test) = 9.23
MSE (Train) = 0.01
MSE (Test) = 0.02
RMSE (Train) = 0.10
RMSE (Test) = 0.15
```

And then, calculating the new top 5 predictors using betas -



# **New Top 5 Predictors**

- 1. 2ndFlrSF
- 2. 1stFlrSF
- 3. MSZoning\_RL
- 4. OverallQual
- 5. MSZoning\_RM

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 -

- 1. A model is **robust** when any variation in the data does not affect its performance in a big way/by a huge margin.
- 2. A **generalisable** model is able to adapt to new, previously unseen data, drawn from the same distribution as the one used to create the model i.e. the variance between the train vs test data scores will not be too much.
- 3. A model should not overfit. If its not overfitting then it can be called robust and generalisable. This is because an overfitting model has very high variance and a smallest change in data affects the model prediction heavily. Such a model will identify all the patterns of a training data, but fail to pick up the patterns in unseen test data.
- 4. The model should not be too complex in order to be robust and generalisable.
- 5. From **Accuracy standpoint**, an overly complex model will have a very high accuracy. So, to make our model more robust and generalisable, we will have to decrease variance which will lead to some bias. Addition of bias means that accuracy will decrease.
- 6. In general, we have to find strike some balance between model accuracy and complexity. This can be achieved by Regularisation techniques like Ridge Regression and Lasso.