

Advanced Regression - Problem Statement

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 to double the value of alpha for both ridge and lasso? What will be the most important predictor variables after the change is implemented?

Answer:

In the case of ridge regression: - When we plot the curve between negative mean absolute error and alpha, we see that as the value of alpha increase from 0 the error term decrease and the train error is showing increasing trend when value of alpha increases when the value of alpha is 2 the test error is minimum so we decided to go with value of alpha equal to 2 for our ridge regression.

For lasso regression I have decided to keep very small value that is 0.01, when we increase the value of alpha the model try to penalize more and try to make most of the coefficient value zero. Initially it came as 0.4 in negative mean absolute error and alpha.

When we double the value of alpha for our ridge regression, 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 dataset 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 r^2 square also decreases.

The most important variable after the changes has been implemented for ridge regression are as follows: -

- 1.MSZoning_FV
- 2.MSZoning_RL
- 3.Neighborhood_Crawfor
- 4.MSZoning_RH

- 5.MSZoning_RM
- 6.SaleCondition_Partial
- 7.Neighborhood_StoneBr
- 8.GrLivArea
- 9.SaleCondition_Normal
- 10.Exterior1st_BrkFace

The most important variable after the changes has been implemented for lasso regression are as follows: -

- 1.GrLivArea
- 2.OverallQual
- 2.OverallCond
- 3.TotalBsmtSF
- 4.BsmtFinSF1
- 5.GarageArea
- 6.Fireplaces
- 7.LotArea
- 8.LotArea
- 9.LotFrontages

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:

Analysis / Observation:

Though the model performance by Ridge Regression was better in terms of R^2 values of Train and Test, it is better to use Lasso, since it brings and assigns a zero value to insignificant features, enabling us to choose the predictive variables. It is always advisable to use simple yet robust model.

It is important to regularize coefficients and improve the prediction accuracy also with the decrease in variance and making the model interpretable.

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 of squares should 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. Ridge regression includes all variables in final model unlike Lasso Regression.

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?

Answer:

The most important predictor variables that will be excluded are: -
GrLivArea - Above grade (ground) living area square feet
OverallQual - Rates the overall material and finish of the house
OverallCond - Rates the overall condition of the house
TotalBsmtSF - Total square feet of basement area
GarageArea - Size of garage in square feet

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:

The model should be as simple as possible, though its accuracy will decrease but it will be more robust and generalisable. It can be also understood using the Bias-Variance trade-off. The simpler the model the more the bias but less variance and more generalizable. Its implication in terms of accuracy is that a robust and generalisable model will perform equally well on both training and test data i.e. the accuracy does not change much for training and test data.

Bias: Bias is error in model when the model is weak to learn from the data. High bias means model is unable to learn details in the data. Model performs poor on training and testing data

Variance: Variance is error in model when model tries to over learn from the data. High variance means model performs exceptionally well on training data as it has very well trained on this of data but performs very poor on testing data as it was unseen data for the model.

It is important to have balance in Bias and Variance to avoid overfitting and underfitting of data.

----- END -----