

Question 1 :- what is the optimal value of alpha for ridge & lasso regression?
what will be the change in the model if you double the value for both ridge & lasso? what will be the most important predictor variable after changing the implementation?

Answer :- During our course of implementing the ridge & lasso regression for Spruce housing Augment dataset

Ridge \rightarrow Alpha = 1

Lasso \rightarrow Alpha = 10

when we calculate the metrics like R^2 Score, RM & RMSE for Ridge regression

the R^2 Score on training data has decreased but it has increased on testing data

when we calculate the metrics like R^2 Score, RM, RMSE for Lasso Regression

the R^2 Score of training data has decreased & it has increased on testing data

The below variables are the most important predictor variables

LotArea
OverallQual
OverallCond
YearBuilt
BsmFltFinFl
TotalBsmtFt
GrLivArea
TotRmsAbvGrnd
StreetPave
RoadMapTotal

Predictions are same But the Coefficients of the predictors has changed.

Question 2: - You have determined the optimal values of λ for ridge & lasso regression during the assignment. Now which one will you choose to apply & why?

Answer: - R^2 score of lasso is slightly higher than for the ridge test data. So we will choose lasso regression over the ridge regression. In order to solve the problem better.

Question 3: - After building the model, you realised that the 5 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 5 most important predictor variables. Which are the 5 most important predictor variables now?

Answer: - LotArea, OverallQual, YearBuilt, BldgType, TotalBsmtSF
post dropping the variables, the R^2 score built has ~~considerable~~
considerable increase in on the training & test data.
Now, the 5 most important variables used for prediction are

1stFlrSF
GrLivArea
StreetName
RoofMtl_Vstal
RoofStyle_gabl.

Question 4:- How can you make sure that a model is robust & generalized? what are the implications of the same for the accuracy of the model & why?

Answer:- the model should be generalized so that the test accuracy is not less than the training score. the model should be accurate for the datasets other than the one which were used during training. Too much importance can't be given to outliers in the dataset. So in turn the accuracy predicted by the model is high.

To ensure this is not the case, the outlier's analysis needs to be done & only those which are relevant to the dataset are retained.

the outlier's which doesn't make sense to be removed from the dataset/model.

If the Best model is not robust then it cannot be trusted for Predictive Analysis.