

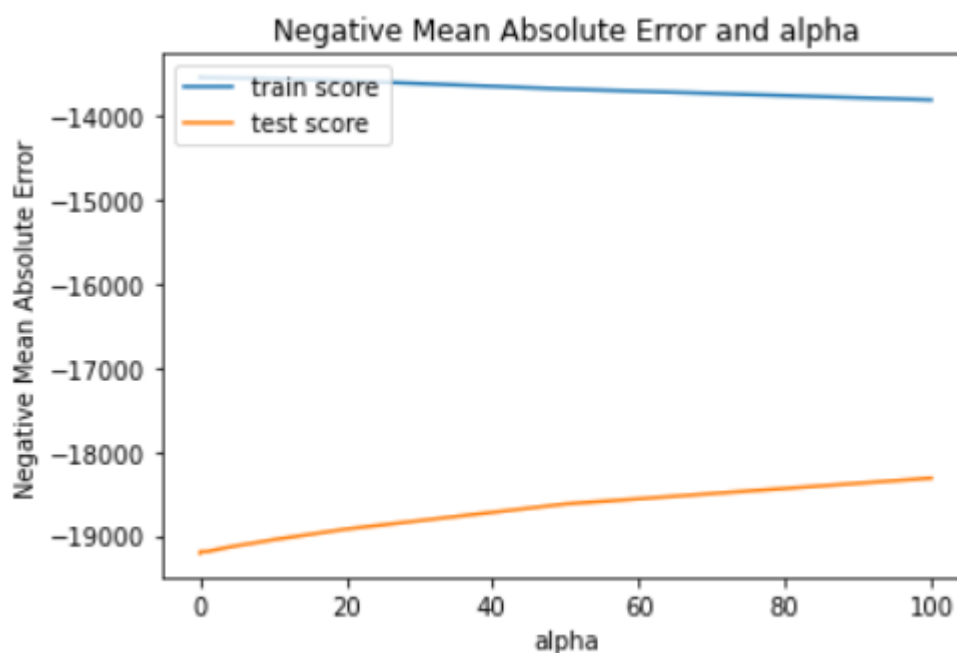
Question 1

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:

1. The optimal value of alpha for ridge and lasso model based on the model built are;
 - a. Lasso: 500
 - b. Ridge: 100
2. There is an increase in both the model but not a very significant level when double the value of alpha
3. The important predictor variables are;

MSSubClass, LotArea, OverallQual, OverallCond, TotalBsmtSF, BsmtFullBath, LowQualFinSF, BsmtHalfBath, YearBuilt_2001-2010, HeatingQC, SaleType, SaleCondition, GarageCond, Exterior1st, Neighborhood, Fireplaces, Condition1, RoofStyle, Exterior2nd_Stone, KitchenQual, Functional_Mod, GarageType_BuiltIn, Heating_Wall, Heating



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:

Lasso involves multiple iterations and puts a penalty over the cost function. Lasso makes the coefficient zero. Hence even though Lasso had a lower r^2 compared to Ridge, I would choose Lasso over ridge. The differences in r^2 between them were not highly significant. Also Lasso works best when there are lots of unwanted variables.

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:

I would go in for the below five variables:

- YearBuilt
- GarageCond
- OverallQual
- HeatingQC
- TotalBsmntSF

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:

Per, Occam's Razor— given two models that show similar 'performance' in the finite training or test data, we should pick the one that makes fewer on the test data due to following reasons: -

- Simpler models are usually more 'generic' and are more widely applicable
- Simpler models require fewer training samples for effective training than the more complex ones and hence are easier to train.
- Simpler models are more robust. Complex models tend to change wildly with changes in the training data set
- Simple models have low variance, high bias and complex models have low bias, high variance
- Simpler models make more errors in the training set. Complex models lead to overfitting — they work very well for the training samples, fail miserably when applied to other test samples

Therefore, to make the model more robust and generalizable, make the model simple but not simpler which will not be of any use.

Regularization can be used to make the model simpler. Regularization helps to strike the delicate balance between keeping the model simple and not making it too naive to be of any use. For regression, regularization involves adding a regularization term to the cost that adds up the absolute values or the squares of the parameters of the model.

Also, Making a model simple lead to Bias-Variance Trade-off:

- A complex model will need to change for every little change in the dataset and hence is very unstable and extremely sensitive to any changes in the training data.
- A simpler model that abstracts out some pattern followed by the data points given is unlikely to change wildly even if more points are added or removed.

Bias quantifies how accurate is the model likely to be on test data. A complex model can do an accurate job prediction provided there is enough training data. Models that are too naïve, for e.g., one that gives same answer to all test inputs and makes no discrimination whatsoever has a very large bias as its expected error across all test inputs are very high.

Variance refers to the degree of changes in the model itself with respect to changes in the training data.

Thus, accuracy of the model can be maintained by keeping the balance between Bias and Variance as it minimizes the total error as shown in the below graph

