

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?

Optimal value of alpha for ridge is 7

Optimal value of alpha for lasso regression is 0.001

After doubling coefficients:

The R2 score came down for both Lasso and Ridge models.

For Ridge:

R2 for train came down from 0.94 to 0.93

R2 for test came down from 0.92 to 0.91

For Lasso:

R2 for train came down from 0.92 to 0.90

R2 for test came down from 0.91 to 0.89

Ridge regression gave the same predictor variables, but Lasso regression removed one predictor and introduced another.

Top predictors are :

OverallQual_9

OverallCond_9

Neighborhood_Crawfor

Functional_Typ

OverallCond_8

GrLivArea

MSZoning_FV

OverallQual_8

Exterior1st_BrkFace

Neighborhood_StoneBr

Question 2

You have determined the optimal value of lambda for ridge and lasso regression during the assignment. Now, which one will you to apply and why?

I used Ridge regression because the R^2 value was slightly higher than that of Lasso regression. Also it helped to reduce magnitude of coefficient.

Another factor is that it brings in bias to the independent variables and hence brings a reliable representation of population

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?

After dropping top 5 predictors, new predictors are :

2ndFlrSF
SaleCondition_Partial
1stFlrSF
OverallCond_9
OverallCond_8

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?

We must be careful not to overfit a model in order to ensure its robustness and generalizability. This is so because an overfitting model has a very large variance and is highly sensitive to even modest changes in the data. Such a model will recognize every pattern in training data, but it will miss the patterns in test data that haven't been observed.

Regularization strategies like Ridge Regression and Lasso could support to provide a tradeoff between model complexity and accuracy.

A model's accuracy will be extremely high if it is overly complex. Therefore, we must reduce variance, which will result in some bias, in order to make our model more reliable and generalizable. Accuracy will decline with the addition of bias.

