

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?

There are no book definition of optimal value, optimal choice of alpha depends on your specific dataset and problem. It's determined through balancing the bias-variance trade-off.

When we double the value of alpha,

With ridge regression, coefficients will be shrunk towards zero, reducing model complexity and potentially improving bias but increasing variance.

With lasso regression, coefficients will be set to exactly zero, resulting in a sparser model with fewer features. This can improve interpretability and reduce overfitting.

W.r.t most important predictor variables,

With ridge regression, all predictors remain in the model, but their importance may be reduced.

With Lasso regression, you can expect even more coefficients to be driven to zero, resulting in a model with fewer predictors.

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?

Ridge regression is suitable when you believe most of the features are relevant and you may want to keep them without excluding entirely. It is strictly convex if this line segment strictly lies above the curve compared to Lasso, so it depends on your optimization requirements. It also add some bias to the model but reduces variance

Lasso regression is suitable if you believe only subset of features are sufficient, resulting in sparser values only with most important predictors. It creates lower-bias model but with potentially higher variance.

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?

Train a new Lasso regression model on the dataset excluding the five most important predictor variables identified in the previous model and examine the magnitude of the coefficients for the remaining predictor variables based on the rank and validate the performance of the new model using cross-validation.

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?

If a model has high accuracy on the training set but low accuracy on the test set, it may be overfitting to the training data. The model has learned specific patterns in the training set that do not generalize well. Ensuring the balance between model complexity and generalizability plays the key part, that can be achieved by understanding the domain, regularization, cross-validation, and thorough evaluation on test sets are essential steps in achieving robust and generalizable models.