

Problem Statement - Part II

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?

Ans: The optimal value of alpha for ridge and lasso regression depends on the dataset being analysed. Alpha controls the strength of the regularization penalty in these models, with higher values of alpha leading to greater regularization and stronger bias-variance trade-offs. To determine the optimal value of alpha, one typically performs cross-validation on a range of alpha values and selects the value that yields the best performance on a validation set. The specific range of alpha values to be tested will depend on the dataset and the desired degree of regularization. If you double the value of alpha for both ridge and lasso, the model will become more regularized, and the coefficients for the predictor variables will shrink further towards zero. This means that the model will be more biased but less prone to overfitting. The most important predictor variables after the change is implemented will depend on the specific dataset and the magnitude of the changes in the coefficients. Generally, the most important predictor variables will be those that have the largest absolute coefficient values, but it is also possible that some predictor variables that were previously unimportant may become more important after the change in regularization strength. It is important to re-evaluate the model performance and interpret the results carefully after any changes are made to the model.

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?

Ans: Ridge regression is typically used when there are many predictor variables that are potentially relevant, and it is expected that most of them will contribute some information to the model. In this case, the goal is to shrink the coefficients towards zero to

reduce overfitting and improve the generalization performance of the model. Ridge regression can also be used to deal with multicollinearity between predictor variables. Lasso regression, on the other hand, is typically used when there are many predictor variables, but it is expected that only a few of them are actually relevant for predicting the outcome variable. In this case, the goal is to select a subset of predictor variables that are most important and exclude the rest. Lasso regression achieves this by shrinking some of the coefficients to exactly zero, effectively removing the corresponding predictor variables from the model. So, depending on the specific situation, I would choose either ridge or lasso regression. If the dataset has many predictor variables, and it is expected that most of them are relevant to the outcome variable, I would choose ridge regression with the optimal value of λ . If the dataset has many predictor variables, but only a few of them are expected to be important, I would choose lasso regression with the optimal value of λ .

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?

Ans: If the five most important predictor variables in the lasso model are not available in the incoming data, we would need to build another model excluding those variables. The most important predictor variables in this new model would depend on the specific dataset and the method used to select the variables. Assuming we are using lasso regression to select the most important variables in the original model, we can fit a new lasso regression model on the data with the five most important predictor variables removed. The most important predictor variables in this new model would be the five variables that have the largest absolute coefficients. It is important to note that the specific variables selected as the most important in the new model may differ from the original model, and the performance of the new model may also be different. It is always a good practice to re-evaluate the model performance and interpret the results carefully after any changes are made to the model.

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?

Ans: To ensure that a model is robust and generalizable, we need to evaluate its performance on data that was not used to train the model. This is typically done using techniques such as cross-validation, where the dataset is split into training and validation sets multiple times, with the model being trained on the training set and evaluated on the validation set in each iteration. This allows us to estimate the generalization performance of the model and detect any potential overfitting to the training set. To further ensure the robustness of the model, we can also use techniques such as bootstrapping or Monte Carlo simulations to estimate the uncertainty of the model estimates and quantify the potential variability in the model performance. The implications of having a robust and generalizable model are significant. A model that is robust and generalizable is less likely to overfit the training data, which can lead to inflated performance estimates and poor generalization to new data. It is important to note that having a highly accurate model on the training data does not necessarily guarantee that the model will perform well on new data. A model that is robust and generalizable is also more likely to be applicable to new datasets and real-world scenarios. This can increase the usefulness and reliability of the model and lead to better decision-making. In summary, ensuring the robustness and generalizability of a model is critical for accurate and reliable predictions, as well as for making informed decisions based on those predictions.