

# Ex4 part (j)

Manuel Alejandro Garcia Acosta

9/10/2020

## Part (j)

Now add some noise to the response. That is, replace **Salary** by **Salary** +  $\epsilon$  where  $\epsilon \sim N(0, \sigma^2)$  for some chosen value of  $\sigma$ . Again, try a standard linear model, ridge, lasso and elastic net where each model is optimally tuned. Repeat this for several increasing values of  $\sigma^2$ . On one plot, plot the error of each model against the amount of additional noise. Does one model or a certain subset of models tend to perform better or worse for large or small values of  $\sigma$ ? Is there some intuition for this? Discuss.

For part (j) I decided to not scale 'Salary'. I just scaled the features and then I added the noise to 'Salary'. Noise was added as  $\epsilon \sim N(0, \sigma^2)$  to 'Salary' for several values of  $\sigma$ . I made two plots of  $\sigma$  vs  $MSE$ , one for small values of  $\sigma$  and the other one for big values.

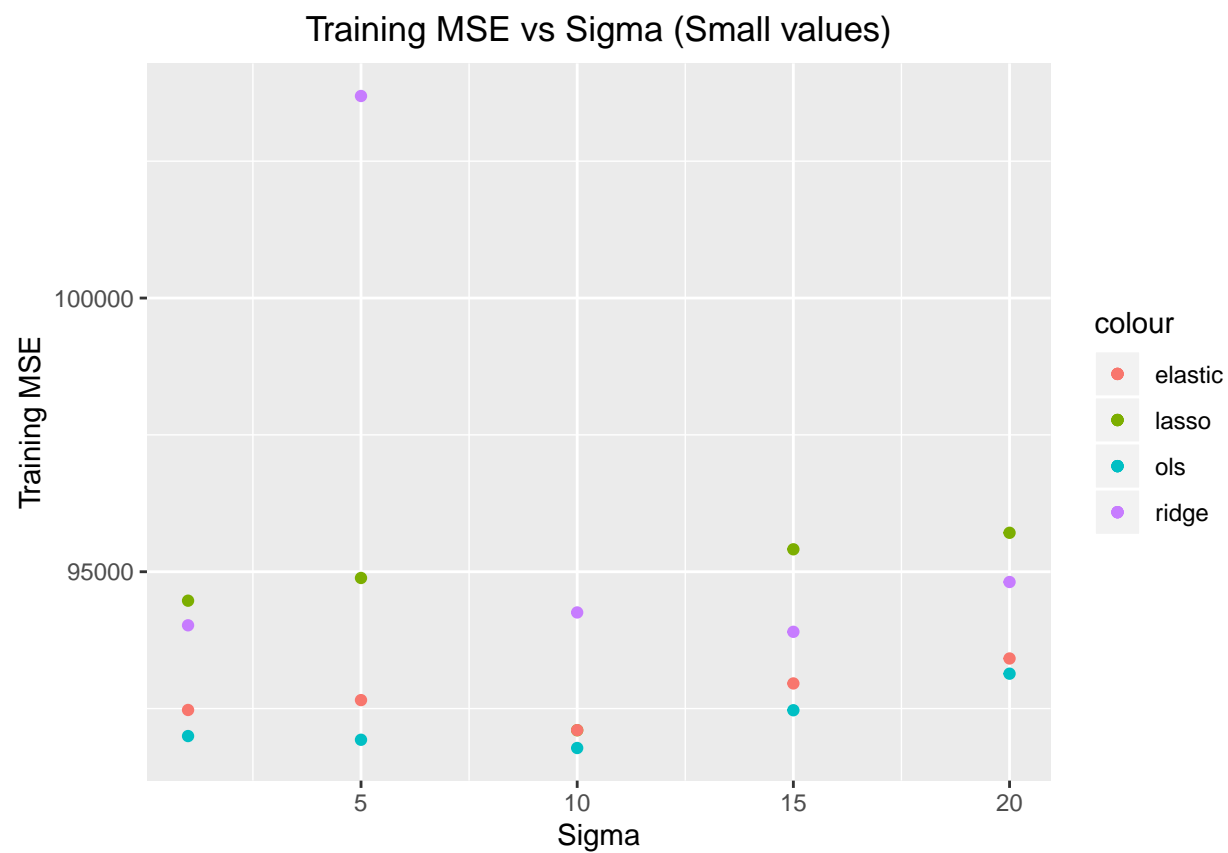
Interestingly, the results were similar to what I got in part (g). OLS and elastic net performed better than ridge regression and lasso for both small and large values of  $\sigma$ . Moreover, for large values of  $\sigma$  all regularization methods start performing way worse than OLS, even elastic net.

After our discussions now I get this happened because I used *training error* to compare the models. Using *test error* *OLS* should perform better for smaller values of  $\sigma$  and the regularization methods should perform better than *OLS* for large values of  $\sigma$ .

### Notes:

- Here I'm using quite a bit of code here, but I'm essentially using the same functions as in parts (b)-(i). I'm omitting the code here but you can find it on the Rmd file 'part\_j.Rmd' I'm attaching.
- For this part I ran all the models using all the original features, i.e., the models were the same as parts (b) and (d)-(f) with additional noise.
- I continued using seeds for reproducibility.

This is the plot for small values of  $\sigma$ , selected values were 1, 5, 10, 15, 20.



This is the plot for large values of  $\sigma$ , selected values were 25, 30, 35, 40, 45.

