

Introduction to Machine Learning - Lab 2

Andrea Bruzzone, Araya Eamrurksiri, Oscar Pettersson, Gustav Sternelöv

Tuesday, November 03, 2015

Assignment 1

1.1

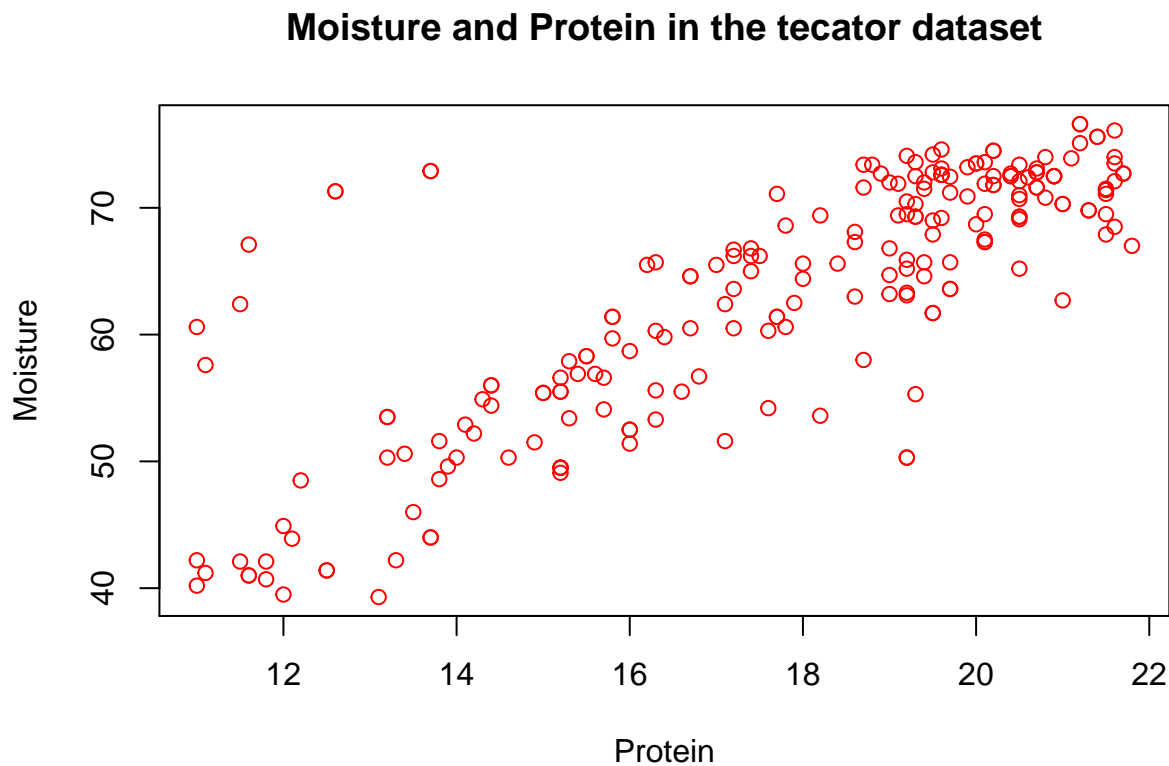
The function `CV()` is created in R. It takes `X`, `Y`, `Lambda` and `Nfolds` as arguments and performs ridge regression by means of `n`-fold cross-validation. `X` is a matrix of predictors, `Y` is a vector of responses, `Lambda` is the shrinkage factor λ and `Nfolds` is `n`. `CV()` returns a score value that is the sum of squared errors.

1.2

Assignment 2

2.1

Following is a plot of Moisture versus Protein:



It can be seen that even though there are some outliers in the beginning of the plot, this data might be well describe by a linear regression.

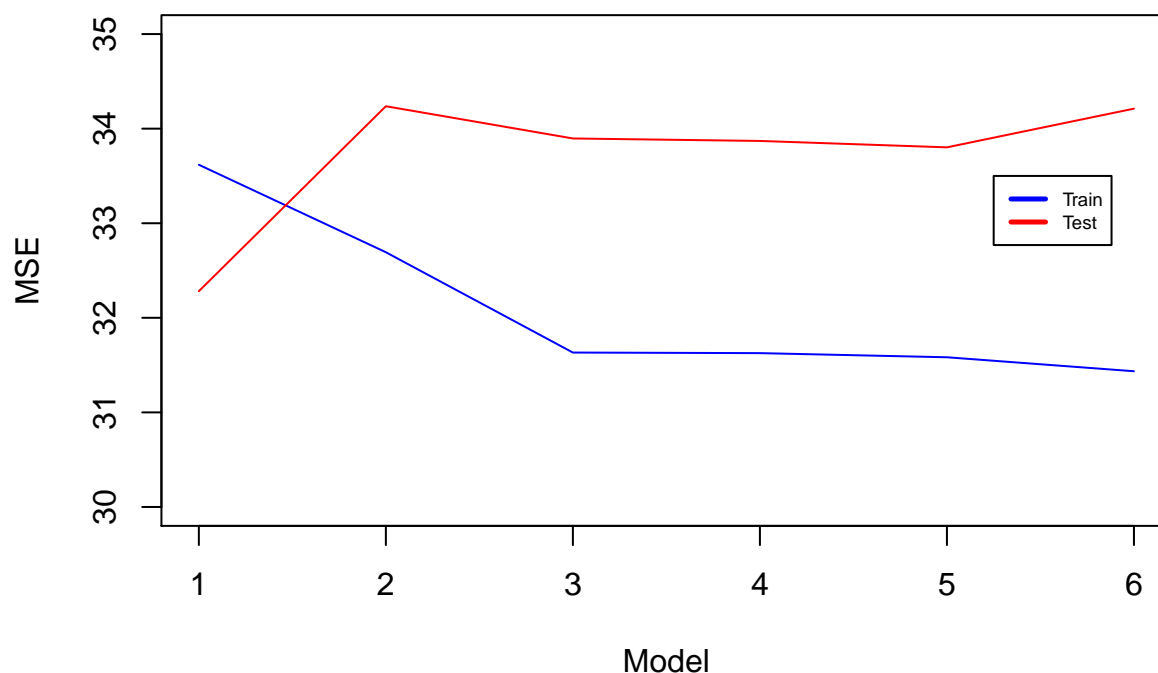
2.2

We consider a model M_i in which expected Moisture is a polynomial function of Protein, including the polynomial terms up to power i .

This model can be described by a normal model as:

$$\mathcal{N}(w_0 + \sum_i w_i x^i, \sigma^2)$$

2.3



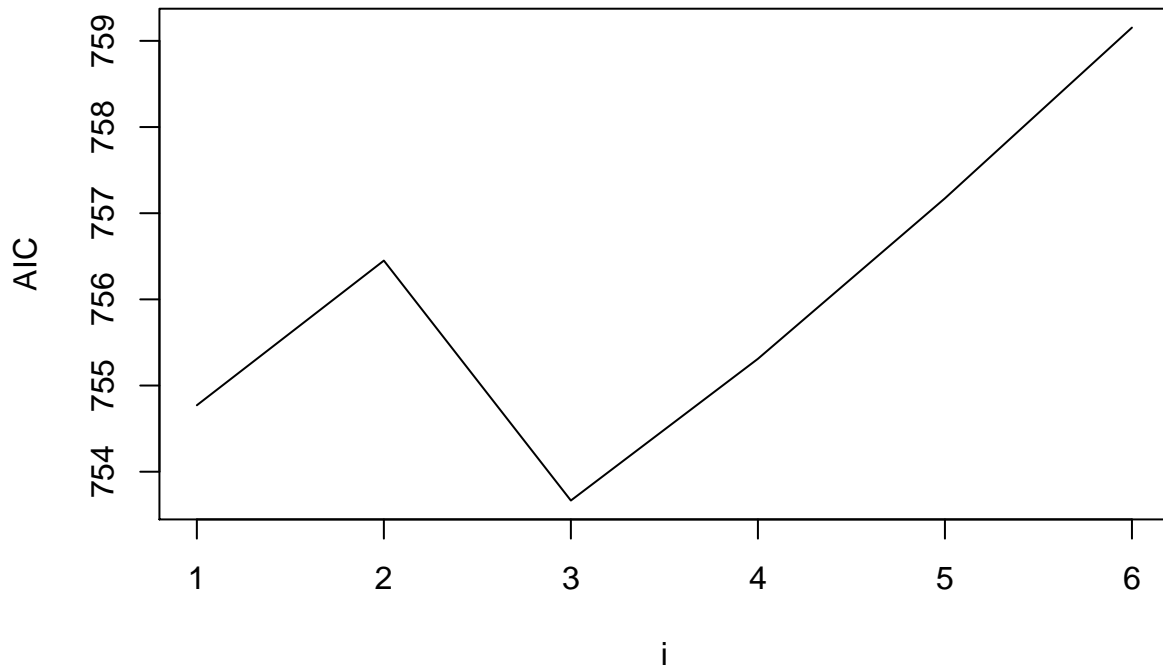
The mean square error for the train dataset appears to be high in the simple model and gradually decreases when the model gets more complex. This is due to the fact that with higher model complexity the model will fit itself closely to the train dataset. Using the model which is too overfitting to train data will result in higher mean square error for the test dataset as can be seen in the plot. Therefore, the model that gives the optimal value of mean square error for both train and test dataset needs to be found.

When consider the bias-variance trade off, this plot suggested that model M_1 , or linear model, is the best among all other models. This model has the lowest value of mean square error for the test dataset. Even though this model has the highest in mean square error in train dataset, it is still better than choosing other models that will overfit data.

2.4

Now the whole data set is used to construct the same models. AIC values is computed to evaluate the models and decide which model that is the best one. The criterion used for comparing the models is that the AIC

value should be as low as possible. With that criterion in mind it is concluded that the model with polynomial terms up to a power of 3 is the preferred model.



2.5

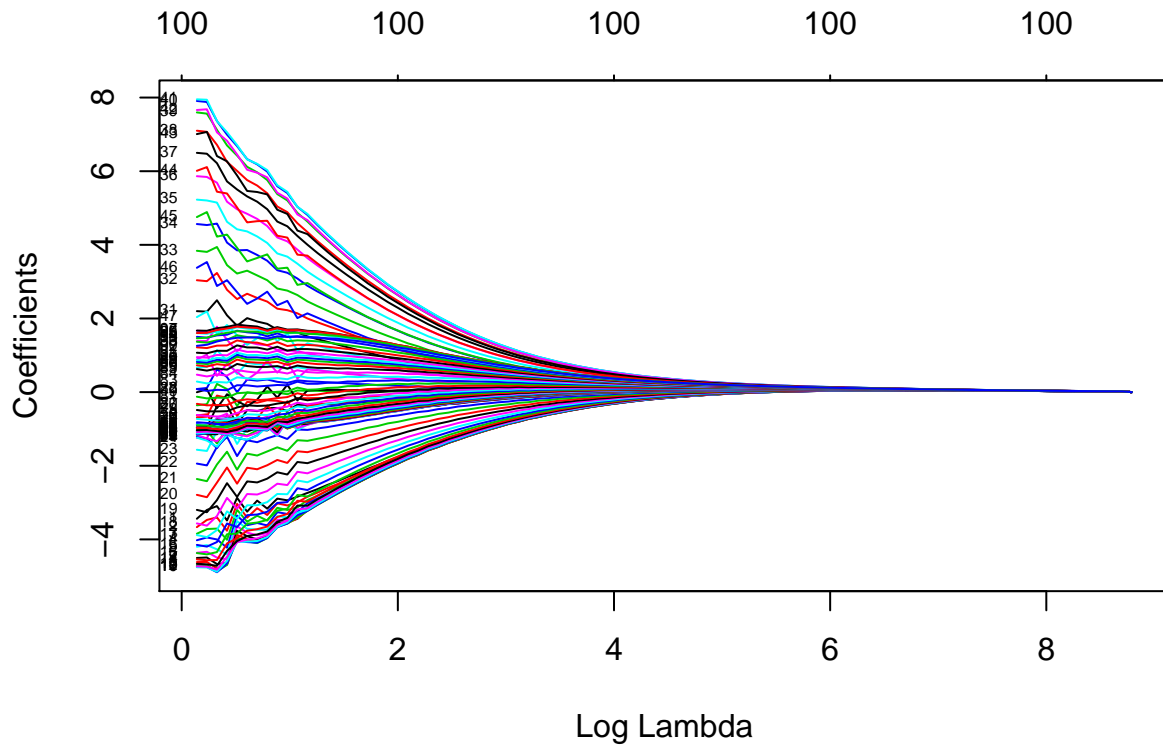
Variable selection is performed for a linear model using `stepAIC` function. `Fat` is response while `Channel11-Channel100` are predictors.

The result of `stepAIC` indicated that 63 variables from 100 variables were selected for this model.

2.6

With the same response and predictor variables as in 2.5 a ridge regression model is fitted. The difference between a ridge regression model and a linear model is the use of a penalty factor λ in the former model. The value of λ affects the coefficients in such a way that they shrink.

How the values of the coefficients in the ridge regression model depends on the log of λ is illustrated by plotting these values.

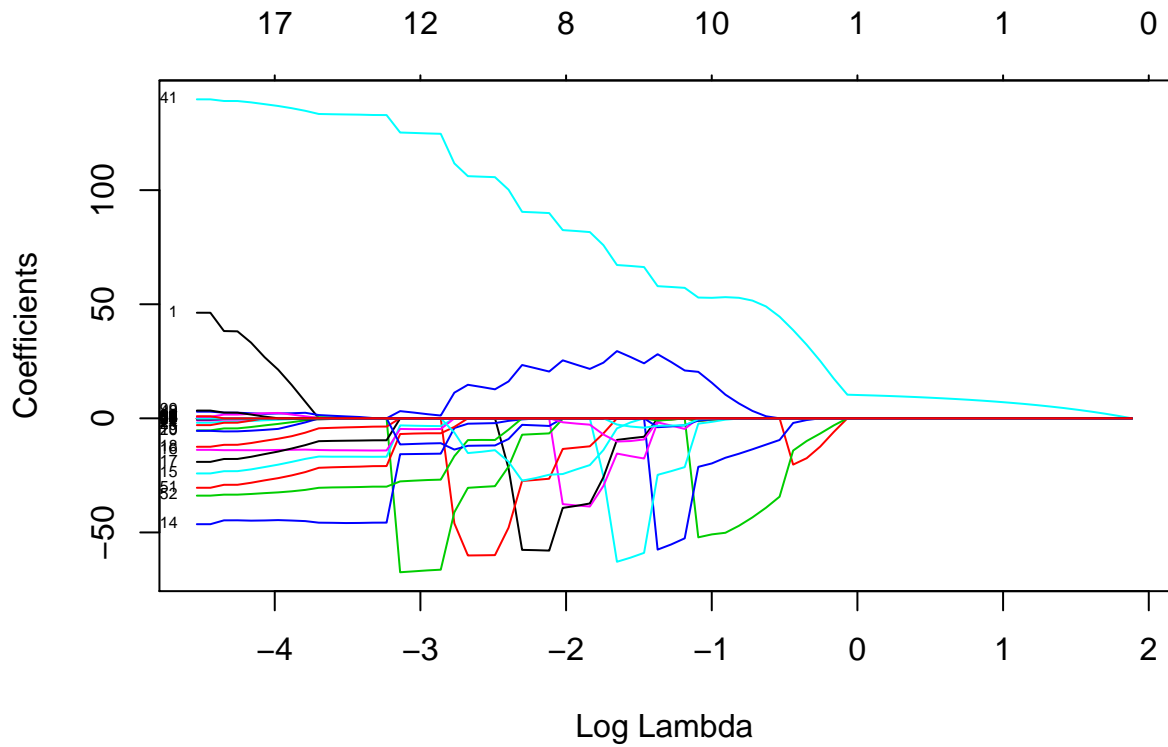


From ridge regression plot, all the coefficients are essentially zero when $\log(\lambda)$ is 8. Ridge regression will always include all the variables in the model, which is 100 variables, and the value of λ selected is indicated by the vertical lines. The plot shows the whole path of variables as they shrink towards zero.

2.7

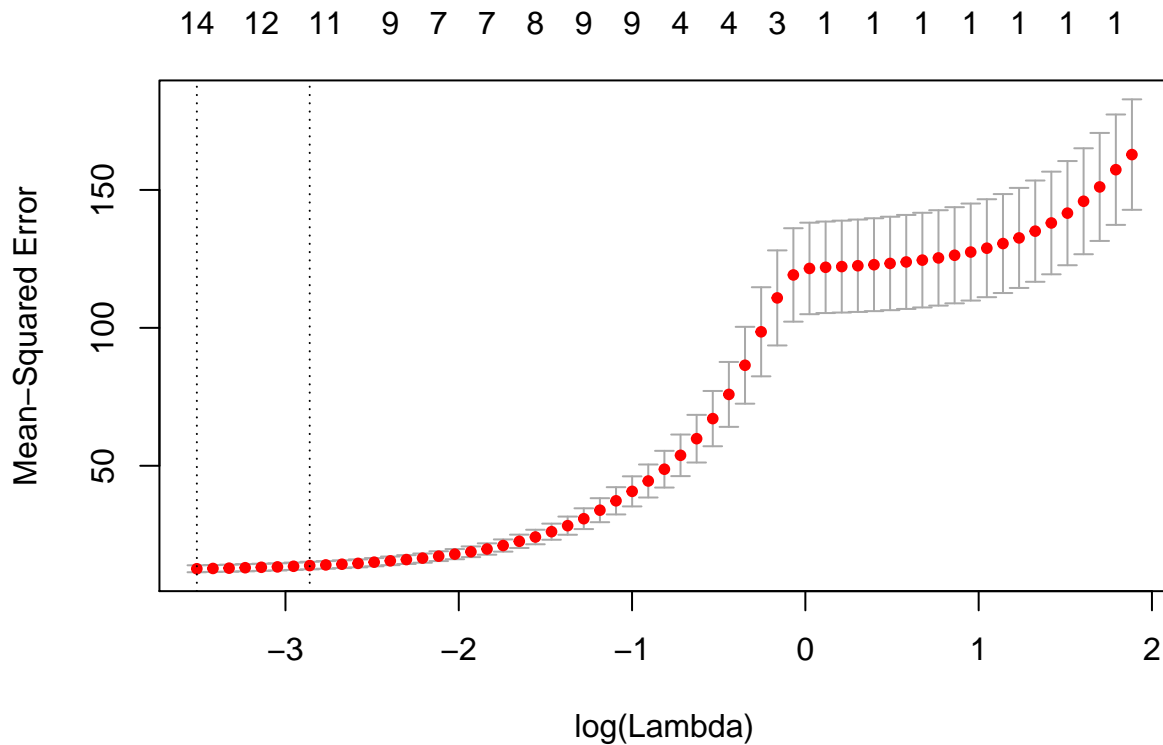
We do the same as in 2.6 but using LASSO instead of the Ridge regression. The idea of the LASSO is very similar to the Ridge regression.

It is used the same function as before, with $\alpha = 1$, that means we use LASSO.



When comparing the plot for the ridge regression and for LASSO the characteristics for the respective methods becomes clear. For the ridge model the coefficients shrinks with higher values of lambda and slowly approaches zero. For the LASSO models coefficients it is harder to see a clear pattern in the dependence between the coefficients and log of λ . The majority of the coefficients quickly equals zero, some has a sort of wiggly curve and some of the coefficients decreases slowly for higher values of log λ .

2.8 Cross-validation is performed to find the optimal LASSO model that is to choose the appropriate lambda and coefficients.



The optimal λ which gives the minimum mean cross-validated error is equal to 0.02985605. Fourteen variables were chosen for an optimal LASSO model.

The plot shows the dependence of the CV score on the log of lambda. The optimal value of λ in terms of $\log(\lambda)$ is -3.511368. When looking at the plot within the one standard error area, it is pretty hard to indicate the exact value of optimal $\log(\lambda)$. It's pretty flat in between. By far, we see that the suggested variables for the optimal model can be varied between 11 to 14 variables with only a slightly difference in CV score. It might also be possible to use other lambda for the model and get the small CV score as the suggested lambda as well. Therefore, it is not quite reliable for the choice of lambda given above.

2.9

Forward-and-backward stepwise regression and cross-validation with LASSO give a different number of variables included in the model. Stepwise regression is a greedy algorithm so it might not find the most optimal number of variables. We also see that ridge regression and LASSO differ. The parameters shrink pretty smoothly for ridge regression while they do so in a more step-wise manner, and eventually become zero, for LASSO.

Group members contribution

Appendix - R-code