

# Assignment 1

FMAN45 - Machine Learning

AUTHOR:

**Dahlberg, Vilmer**

`vi8808da-s@student.lu.se`

April 20, 2021

## Introduction

In this report the LASSO method is studied, and is used to solve two regression problems. Firstly, a set of noisy data consisting of two main frequencies is to be denoised. The hyperparameter used in LASSO is first chosen using a guess-and-check method, but later trained using a cross-validation scheme. Next, a clip of music is denoised using LASSO, where the hyperparameter is trained using the same cross-validation scheme. The report consists of the seven tasks stated in the assignment, which are answered and discussed one by one.

## Task 1

The LASSO problem can be solved using a coordinate-wise iterative method, where the  $i$ :th parameter is given by the solution to

$$\underset{w_i}{\text{minimize}} \frac{1}{2} \|\mathbf{r}_i - \mathbf{x}_i w_i\|_2^2 + \lambda |w_i|. \quad (1)$$

where  $\mathbf{r}_i = \mathbf{t} - \sum_{l \neq i} \mathbf{x}_l w_l$  and  $\mathbf{x}_i$  is the  $i$ :th row of the regression matrix  $\mathbf{X}$ .

## Solution

Assume the  $j$ :th iterate of the  $i$ :th parameter is sought, meaning  $w_i^j$ . Let  $\mathbf{r}^{(j-1)} = \mathbf{t} + \sum_{l < i} \mathbf{x}_l w_l^j + \sum_{l > i} \mathbf{x}_l w_l^{(j-1)}$  be the residual at the current step. Then, the next iterate is found by minimizing the function below with respect to  $w_i^j$

$$f(w_i^j) = \frac{1}{2} \|\mathbf{r}_i^{(j-1)} - \mathbf{x}_i w_i^j\|_2^2 + \lambda |w_i^j| = \frac{1}{2} (\mathbf{r}_i^{(j-1)} - \mathbf{x}_i w_i^j)^T (\mathbf{r}_i^{(j-1)} - \mathbf{x}_i w_i^j) + \lambda |w_i^j|$$

As this is a convex problem, the minimum exists and is found by setting the derivative to zero.

Assuming  $w_i^j \neq 0 \iff \frac{d|w_i^j|}{dw_i^j} = \frac{w_i^j}{|w_i^j|}$ , we get

$$\begin{aligned} \frac{df}{dw_i^j} &= -\mathbf{x}_i^T \mathbf{r}_i^{(j-1)} + \mathbf{x}_i^T \mathbf{x}_i w_i^j + \lambda \frac{w_i^j}{|w_i^j|} = 0 \iff \\ w_i^j (\mathbf{x}_i^T \mathbf{x}_i + \lambda \frac{1}{|w_i^j|}) &= \mathbf{x}_i^T \mathbf{r}_i^{(j-1)} \end{aligned} \quad (2)$$

To get an explicit expression for  $w_i^j$ , move the first term to the right-hand side and take the absolute value of the whole expression.

$$|w_i^j| \left| \mathbf{x}_i^T \mathbf{x}_i + \lambda \frac{1}{|w_i^j|} \right| = |\mathbf{x}_i^T \mathbf{r}_i^{(j-1)}| \quad (3)$$

Consider the left hand side of the expression. Using the fact that  $\mathbf{x}_i^T \mathbf{x}_i, \lambda > 0$  we find

$$|w_i^j| \left| \mathbf{x}_i^T \mathbf{x}_i + \lambda \frac{1}{|w_i^j|} \right| = |\mathbf{x}_i^T \mathbf{x}_i| |w_i^j| + \lambda = \mathbf{x}_i^T \mathbf{x}_i |w_i^j| + \lambda$$

Finally solve for  $|w_i^j|$  in equation 3,  $|w_i^j| = \frac{1}{\mathbf{x}_i^T \mathbf{x}_i} (|\mathbf{x}_i^T \mathbf{r}_i^{(j-1)}| - \lambda)$ , and insert into equation 2 to get

$$\begin{aligned}
 w_i^j (\mathbf{x}_i^T \mathbf{x}_i + \lambda \frac{1}{|w_i^j|}) &= \mathbf{x}_i^T \mathbf{r}_i^{(j-1)} && \iff \\
 w_i^j (\mathbf{x}_i^T \mathbf{x}_i |w_i^j| + \lambda) &= |w_i^j| \mathbf{x}_i^T \mathbf{r}_i^{(j-1)} && \iff \\
 w_i^j |\mathbf{x}_i^T \mathbf{r}_i^{(j-1)}| &= \frac{1}{\mathbf{x}_i^T \mathbf{x}_i} (|\mathbf{x}_i^T \mathbf{r}_i^{(j-1)}| - \lambda) \mathbf{x}_i^T \mathbf{r}_i^{(j-1)} && \iff \\
 w_i^j &= \frac{\mathbf{x}_i^T \mathbf{r}_i^{(j-1)}}{\mathbf{x}_i^T \mathbf{x}_i |\mathbf{x}_i^T \mathbf{r}_i^{(j-1)}|} (|\mathbf{x}_i^T \mathbf{r}_i^{(j-1)}| - \lambda)
 \end{aligned}$$

## Task 2

Consider the case where the regression matrix is an orthonormal basis, i.e.  $\mathbf{X}$  is square and  $\mathbf{X}^T \mathbf{X} = \mathbf{I}$ . Show that the coordinate descent will converge in at most 1 iteration.

### Solution

Begin by considering the update formula

$$w_i^j = \begin{cases} \frac{\mathbf{x}_i^T \mathbf{r}_i^{(j-1)}}{\mathbf{x}_i^T \mathbf{x}_i |\mathbf{x}_i^T \mathbf{r}_i^{(j-1)}|} (|\mathbf{x}_i^T \mathbf{r}_i^{(j-1)}| - \lambda), & |\mathbf{x}_i^T \mathbf{r}_i^{(j-1)}| > \lambda \\ 0 & \text{else} \end{cases}$$

where  $\mathbf{r}^{(j-1)} = \mathbf{t} + \sum_{l < i} \mathbf{x}_l w_l^j + \sum_{l > i} \mathbf{x}_l w_l^{(j-1)}$ . Using the fact that  $\mathbf{X}$  is an orthonormal basis, meaning  $\mathbf{x}_j^T \mathbf{x}_k = \delta_{jk}$ , compute the product in the numerator.

$$\mathbf{x}_i^T \mathbf{r}_i^{(j-1)} = \mathbf{x}_i^T \mathbf{t} + \sum_{l < i} \mathbf{x}_i^T \mathbf{x}_l w_l^j + \sum_{l > i} \mathbf{x}_i^T \mathbf{x}_l w_l^{(j-1)} = \mathbf{x}_i^T \mathbf{t}.$$

Thus, the update formula becomes

$$w_i^j = \begin{cases} \frac{\mathbf{x}_i^T \mathbf{t}}{|\mathbf{x}_i^T \mathbf{t}|} (|\mathbf{x}_i^T \mathbf{t}| - \lambda), & |\mathbf{x}_i^T \mathbf{t}| > \lambda \\ 0 & \text{else.} \end{cases}$$

Which is obviously independent on the iteration number  $j$ , and hence the iteration converges in at most 1 iteration.

## Task 3

Show that the LASSO estimate's bias, for an orthogonal regression matrix and data generated by  $\mathbf{t} = \mathbf{X}\mathbf{w}^* + \mathbf{e}$ ,  $\mathbf{e} \sim \mathcal{N}(\mathbf{0}_N, \sigma \mathbf{I}_N)$ , is given by

$$\lim_{\sigma \rightarrow 0} E \left( \hat{w}_i^1 - w_i^* \right) = \begin{cases} -\lambda, & w_i^* > \lambda \\ -w_i^*, & |w_i^*| \leq \lambda \\ \lambda, & w_i^* < -\lambda \end{cases} \quad \forall i. \quad (4)$$

### Solution

As shown above when the regression matrix is orthogonal the update formula for the parameters  $w_i^j$  is given by

$$w_i^j = \begin{cases} \frac{\mathbf{x}_i^T \mathbf{t}}{|\mathbf{x}_i^T \mathbf{t}|} (|\mathbf{x}_i^T \mathbf{t}| - \lambda), & |\mathbf{x}_i^T \mathbf{t}| > \lambda \\ 0 & \text{else.} \end{cases}$$

Begin by considering the case where  $|\mathbf{x}_i^T \mathbf{t}| \leq \lambda$ , then

$$\lim_{\sigma \rightarrow 0} E(\hat{w}_i^1 - w_i^*) = E(-w_i^*) = -w_i^*.$$

Next consider the case where  $|\mathbf{x}_i^T \mathbf{t}| > \lambda$ . The model used to generate  $\mathbf{t}$  is known, so the product  $\mathbf{x}_i^T \mathbf{t}$  can easily be computed

$$\mathbf{x}_i^T \mathbf{t} = \mathbf{x}_i^T (\mathbf{X} \mathbf{w}^* + \mathbf{e}) = \sum_j \mathbf{x}_i^T \mathbf{x}_j w_j + \mathbf{x}_i^T \mathbf{e} = w_i^* + \mathbf{x}_i^T \mathbf{e},$$

where it was used that  $\mathbf{x}_i^T \mathbf{x}_j = \delta_{ij}$ . Now, study the limit of the bias as the variation in the noise tends to zero,

$$\lim_{\sigma \rightarrow 0} E(\hat{w}_i^1 - w_i^*) = \lim_{\sigma \rightarrow 0} E\left(\frac{w_i^* + \mathbf{x}_i^T \mathbf{e}}{|\mathbf{x}_i^T \mathbf{t}|} (|\mathbf{x}_i^T \mathbf{t}| - \lambda) - w_i^*\right).$$

Since the noise  $\mathbf{e}$  has the expected value  $\mathbf{0}$ , the product  $\mathbf{x}_i^T \mathbf{e}$  tends to zero as the variation in the noise tends to zero, implying  $\mathbf{x}_i^T \mathbf{t}$  tends to  $w_i^*$ , thus

$$\lim_{\sigma \rightarrow 0} E(\hat{w}_i^1 - w_i^*) = E\left(\frac{w_i^*}{|w_i^*|} (|w_i^*| - \lambda) - w_i^*\right) = E\left(w_i^* \left(1 - \frac{\lambda}{|w_i^*|}\right) - w_i^*\right) = E\left(-w_i^* \frac{\lambda}{|w_i^*|}\right) = \text{sign}(-w_i^*) \lambda.$$

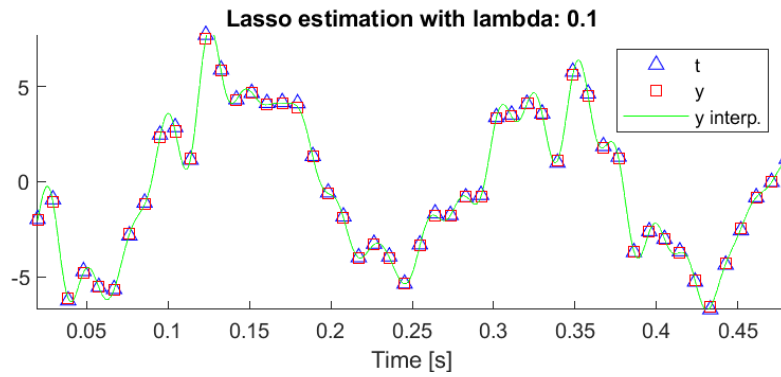
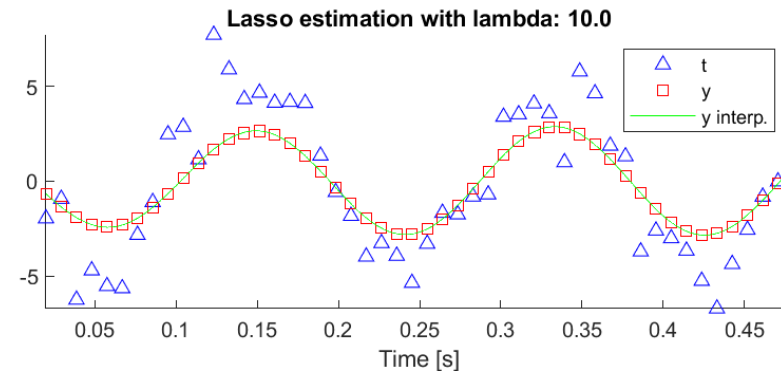
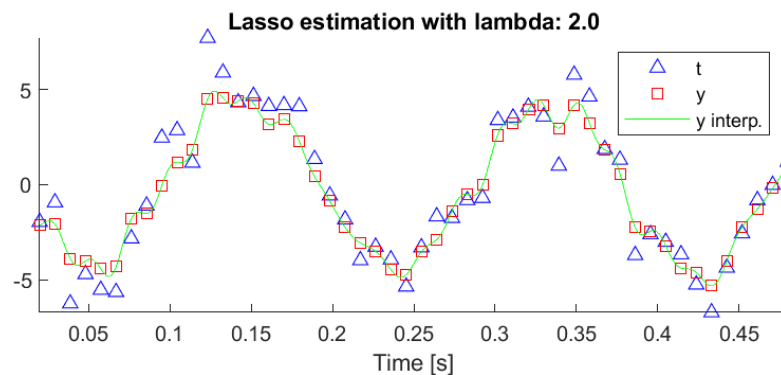
Divide the case  $|\mathbf{x}_i^T \mathbf{t}| > \lambda$  into the two cases  $\mathbf{x}_i^T \mathbf{t} > \lambda$ , or  $\mathbf{x}_i^T \mathbf{t} < -\lambda$ , and use the fact that  $\lambda$  is positive.

$$\begin{cases} w_i^* > \lambda \implies w_i > 0 \implies \lim_{\sigma \rightarrow 0} E(\hat{w}_i^1 - w_i^*) = \text{sign}(-w_i^*) \lambda = -\lambda \\ w_i^* < -\lambda \implies w_i < 0 \implies \lim_{\sigma \rightarrow 0} E(\hat{w}_i^1 - w_i^*) = \text{sign}(-w_i^*) \lambda = \lambda \end{cases}$$

In words, the absolute value of the bias tends to  $\lambda$  if the absolute value of the weight is larger than  $\lambda$ , otherwise the estimation bias is the weight itself. This seems reasonable as the acronym LASSO stands for least absolute shrinkage and selection operator.

## Task 4

1. The LASSO estimation was implemented and tested using the given data. The results are shown below.

(a) Small penalty  $\lambda = 0.1$ .(b) A large penalty  $\lambda = 10$ .(c) A slightly smaller penalty  $\lambda = 2$ .Figure 1: LASSO estimation using different  $\lambda$  values.

The penalty coefficient  $\lambda$  determines how to weigh the error in the estimation and the magnitude of the coefficients  $w$ . If a small  $\lambda$  is used, minimizing the error will be prioritized over minimizing the magnitude of the weights, and vice versa.

From the first plot, figure 1a, the reconstruction follows the data very closely. However, using this estimation for prediction would give poor results, since the data is overfit to include the random noise. A larger penalty should be used to decrease the amount of frequencies included.

In the next plot, figure 1b, the reconstruction seems to follow the low frequency oscillation well, but the high frequency oscillation is not followed. The reconstruction does not represent the model, and the error is too large. A smaller penalty should be used.

In the final plot, figure 1c, the reconstruction seems to follow both the low and high frequency oscillations well, although there is still some error. This  $\lambda$  balances the minimization of the

error and the minimization of the magnitude of the coefficients well.

2. The number of nonzero coefficients were counted for each lambda value.

lambda	num. of $w \neq 0$
0.1	263
10	5
2	44

The number of coefficients needed to model the data is four, two for each frequency. However, around 44 coefficients are needed to get a good reconstruction of the data using this model. This also highlights the difficulty in choosing the penalty  $\lambda$ .

## Task 5

1. The figure below shows the root-mean-square error for the estimation folds in red and the validation folds in blue for different values of lambda. The optimal lambda is chosen as the lambda that minimizes the validation error.

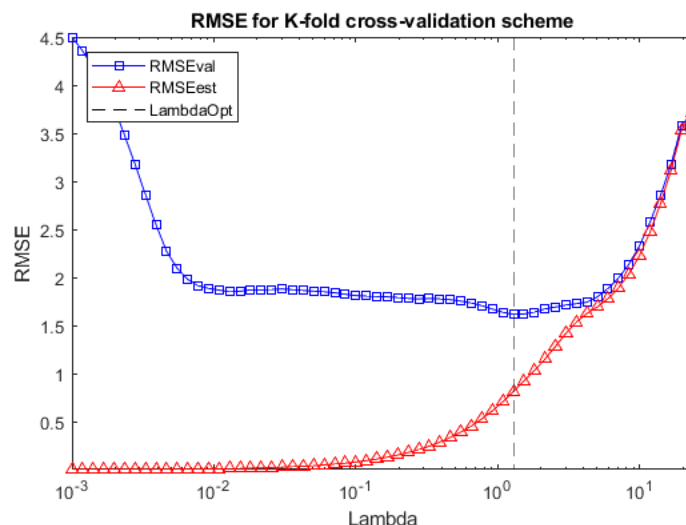
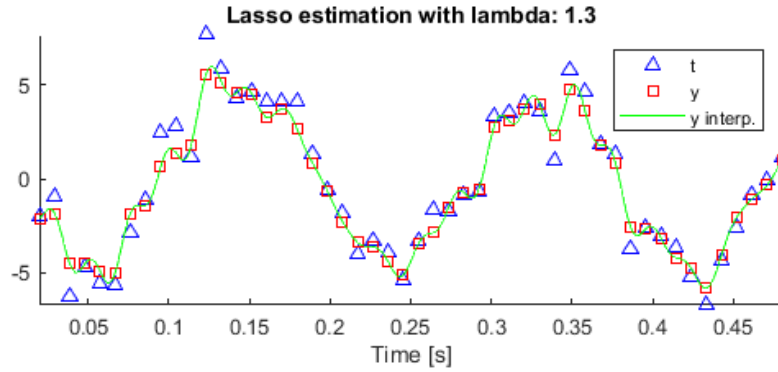


Figure 2: RMSE for validation and estimation folds. The lambda which minimizes  $RMSE_{val}$  is chosen as the optimal lambda.

The estimation error grows as lambda increases, which is expected since an increase in lambda penalizes the magnitude of the coefficients and sacrifices the error. The validation error changes very little at and near the minimum, and choosing a lower lambda value would have a similar validation error but could reduce the estimation value by a significant amount.

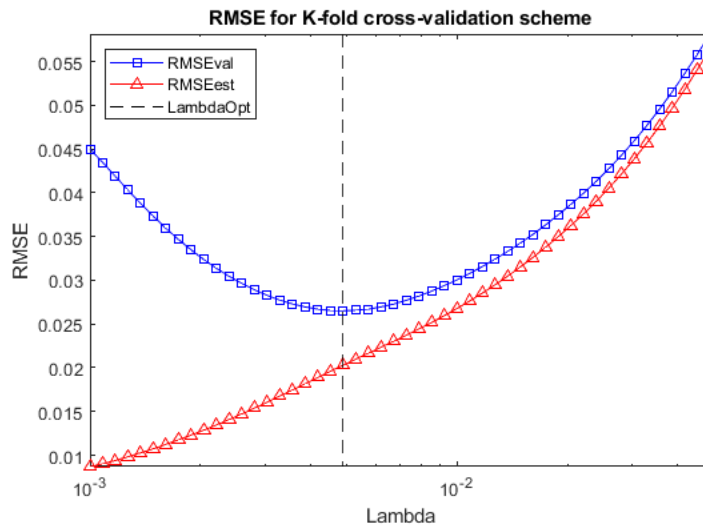
Also note that the estimation error is always smaller than the validation error, which is especially clear when lambda is small. This is because a small lambda, hence a small penalisation, tends to result in an overfit of the data. If the data is overfit the estimation error shrinks, but the validation error grows since the noise is random. If the lambda is very small, the estimation error becomes very small, but the validation error grows. This is again due to the fact that the noise is random and an overfit of the noise in the estimation fold will result in large errors in the validation folds.

2. The plot shows a reconstruction of the data using the optimal lambda value. The number of nonzero  $w$  is still far too large. This suggests that the problem might lie in the model, and not the choice of lambda.

Figure 3: Number of nonzero  $w$ : 70.

## Task 6

Just as in the previous task, the estimation error grows as the penalty increases, while the validation error has a clear minimum. Likewise, the estimation error is always lower than the validation error. Although the minimum is clearer in this case. The optimal lambda was about  $4.9e-3$ .

Figure 4: RMSE for validation and estimation folds. The lambda which minimizes  $RMSE_{val}$  is chosen as the optimal lambda.

## Task 7

A clear difference is heard when the cleaned audio is played after the noisy audio. However, the audio is still a bit noisy. Increasing the penalty reduces the noise, but reduces the amount of frequencies in the reconstruction which results in the piano sounding off. On the other hand, reducing the penalty makes the piano sound better but the noise comes back.

The method works, but has flaws and does not perform very well. A different approach is needed to remove all noise without destroying the piano sound. The method definitely works as a black-box method, but more information about the piano sound can be used. For example, the actual frequencies which comprise the piano sound are well known and can be used to further improve the quality in the cleaning of the audio.