

# Week 3 Homework Submission

## Linear Regression

Iuliia Skobleva  
Matriculation Number 03723809

November 10, 2019

### Least squares regression

#### Problem 1

The sum of squares error function for a dataset, where each point is weighted is given by:

$$E_{weighted}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N t_i [\mathbf{w}^T \phi(\mathbf{x}_i) - y_i]^2 \quad (1)$$

And we want to find  $\mathbf{w} = \arg \min E_{weighted}(\mathbf{w})$ . In matrix form, the equation 1 can be written as:

$$\begin{aligned} E_{weighted}(\mathbf{w}) &= (\mathbf{y} - \Phi \mathbf{w})^T T (\mathbf{y} - \Phi \mathbf{w}) \\ &= (\mathbf{y}^T - \mathbf{w}^T \Phi^T) T (\mathbf{y} - \Phi \mathbf{w}) \\ &= (\mathbf{y}^T - \mathbf{w}^T \Phi^T) (T \mathbf{y} - T \Phi \mathbf{w}) \end{aligned} \quad (2)$$

Using the matrix cookbook, we find the gradient and set it to 0 to find  $\mathbf{w}$ :

$$\nabla_{\mathbf{w}} E_{weighted} = \Phi^T T \Phi \mathbf{w} - \Phi^T T \mathbf{y} \stackrel{!}{=} 0 \quad (3)$$

In the end, we have:

$$\mathbf{w} = (\Phi^T T \Phi)^{-1} \Phi^T T \mathbf{y} \quad (4)$$

Another question is how this weighting factor  $t_i$  can be interpreted in terms of the variance of the noise on the data and data points for which there are exact copies in

the dataset. If we choose the weighting factor appropriately, meaning we give lower weight to data points with high variance, we can derive more accurate predictions,  $\mathbf{y} = \mathbf{X}\mathbf{w}$ . If we have exact copies in the data set, the analysis will be inaccurate, since we cannot differentiate between points that are just noise (with high variance) and points that are a different class. There, applying the weighting might lead to misleading results.

## Ridge Regression

### Problem 2

In this problem we augment  $\Phi \in \mathbb{R}^{N \times M}$  with  $\sqrt{\lambda}\mathbf{I}_{M \times M}$  and the target vector  $\mathbf{y}$  with  $M$  rows of 0s. Now we can show that the ridge regression estimates can be obtained by ordinary least squared regression.

So, using  $E_{RR}(\mathbf{w}) = \frac{1}{2}(\hat{\mathbf{y}} - \hat{\Phi}\mathbf{w})^T(\hat{\mathbf{y}} - \hat{\Phi}\mathbf{w})$ , where new  $\hat{\Phi}$  and  $\hat{\mathbf{y}}$  are given as,  $\hat{\Phi} = \begin{pmatrix} \Phi \\ \sqrt{\lambda}\mathbf{I}_{M \times M} \end{pmatrix}$  and  $\hat{\mathbf{y}} = \begin{pmatrix} \mathbf{y} \\ \mathbf{0}_{M \times 1} \end{pmatrix}$  we find:

$$\begin{aligned} E_{RR}(\mathbf{w}) &= \frac{1}{2} \left[ \begin{pmatrix} \Phi \\ \sqrt{\lambda}\mathbf{I} \end{pmatrix} \mathbf{w} - \begin{pmatrix} \mathbf{y} \\ \mathbf{0} \end{pmatrix} \right]^T \left[ \begin{pmatrix} \Phi \\ \sqrt{\lambda}\mathbf{I} \end{pmatrix} \mathbf{w} - \begin{pmatrix} \mathbf{y} \\ \mathbf{0} \end{pmatrix} \right] \\ &= \frac{1}{2} \left[ \begin{pmatrix} \Phi\mathbf{w} - \mathbf{y} \\ \sqrt{\lambda}\mathbf{I}\mathbf{w} - \mathbf{0} \end{pmatrix} \right]^T \left[ \begin{pmatrix} \Phi\mathbf{w} - \mathbf{y} \\ \sqrt{\lambda}\mathbf{I}\mathbf{w} - \mathbf{0} \end{pmatrix} \right] \\ &= \frac{1}{2} \left[ (\Phi\mathbf{w} - \mathbf{y})^T(\Phi\mathbf{w} - \mathbf{y}) + \lambda\mathbf{w}^T\mathbf{w} \right] \end{aligned} \tag{5}$$

Thus, we find the least squares loss function with  $L_2$ -regularization.

### Problem 3

We now want to derive the closed form solution for the expression in equation 5, aka find:

$$\mathbf{w} = \arg \min E_{RR}(\mathbf{w}) \tag{6}$$

This is done by calculating the gradient of  $E_{RR}(\mathbf{w})$ .

$$\nabla_{\mathbf{w}} E_{RR}(\mathbf{w}) = \nabla_{\mathbf{w}} \frac{1}{2} \left[ (\Phi\mathbf{w} - \mathbf{y})^T(\Phi\mathbf{w} - \mathbf{y}) + \lambda\mathbf{w}^T\mathbf{w} \right] \stackrel{!}{=} 0 \tag{7}$$

We find with the help of the matrix cookbook:

$$\Phi^T \Phi \mathbf{w} - \Phi^T \mathbf{y} + \lambda \mathbf{I} \mathbf{w} \stackrel{!}{=} 0 \quad (8)$$

This gives the solution for the closed form expression:

$$\mathbf{w} = (\Phi^T \Phi + \lambda \mathbf{I})^{-1} \Phi^T \mathbf{y} \quad (9)$$

Another question is what happens when the number of training samples  $N$  is smaller than the number of basis functions  $M$ . In this case, our model is too flexible and this leads to overfitting. To avoid that, we need to specify the regularization strength  $\lambda$  and make it large enough so that the weights  $\mathbf{w}$  are smaller and the model shows small errors on the validation set.

## Comparison of Linear Regression Models

### Problem 5

We assume that we have fitted an  $L_2$ -regularized linear regression model and obtained the optimal weight vector  $\mathbf{w}^*$ :

$$\mathbf{w}^* = \arg \min \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i - y_i)^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \quad (10)$$

- If we created a new matrix  $\mathbf{X}_{\text{new}} = a\mathbf{X}$ , what weight vector would produce the same results on  $\mathbf{X}_{\text{new}}$  as  $\mathbf{w}^*$  on  $\mathbf{X}$ ? To find out, we look at equation 9 (which is the closed form solution for ridge regression error function), where instead of having the function  $\Phi$ , we have the matrix  $\mathbf{X}$  thus getting  $\mathbf{w}^* = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$ . The new weight vector would be:

$$\mathbf{w}_{\text{new}}^* = (a^2 \mathbf{X}^T \mathbf{X} + \lambda_{\text{new}} \mathbf{I})^{-1} \cdot a \mathbf{X}^T \mathbf{y} \quad (11)$$

- The predictions are given by  $\mathbf{y} = \mathbf{X} \mathbf{w}$  and for the new matrix it is  $\mathbf{y}_{\text{new}} = \mathbf{X}_{\text{new}} \mathbf{w}_{\text{new}}$ . The prediction for the old dataset equals:

$$\begin{aligned} \mathbf{y} &= \mathbf{X} \mathbf{w} \\ &= \mathbf{X} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \cdot \mathbf{X}^T \mathbf{y} \end{aligned} \quad (12)$$

The solution for the weight vector  $\mathbf{w}_{\text{new}}$  is found in equation 11. For the prediction we have:

$$\begin{aligned}\mathbf{y}_{\text{new}} &= a\mathbf{X}\mathbf{w}_{\text{new}} \\ &= a\mathbf{X}(a^2\mathbf{X}^T\mathbf{X} + \lambda_{\text{new}}\mathbf{I})^{-1} \cdot a\mathbf{X}^T\mathbf{y}\end{aligned}\tag{13}$$

In the above equation we see that if  $\lambda_{\text{new}} = a^2\lambda$ , then  $\mathbf{y} = \mathbf{y}_{\text{new}}$ .