

3<sup>rd</sup> Homework

Thomas Saltos

## Exercise 1

We want to minimize the  $L(\theta)$ , so we have to set the derivative to 0:

$$\frac{\partial L(\theta)}{\partial \theta} = 0 \Rightarrow 2 \sum_{n=1}^N ((-x_n)(y_n - x_n^T \theta) + \lambda \theta) = 0 \Rightarrow \sum_{n=1}^N (x_n x_n^T + \lambda \theta) = \sum_{n=1}^N x_n y_n$$

Let's define the matrix  $X = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{pmatrix}$ . Now the  $\sum_{n=1}^N x_n x_n^T$  can be written in matrix form as  $X^T X$ .

Moreover, the term  $\sum_{n=1}^N x_n y_n$  equals  $X^T y$  where  $y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}$ .

The ridge-regression solution can be expressed as  $\theta = (X^T X + \lambda I)^{-1} X^T y$

## Exercise 2

- The second term (biased term) equals to 0, since  $\theta_{MVU}$  is unbiased estimator of  $\theta_0$ , thus  $E(\theta_{MVU}) = \theta_0$
- $E(\theta_b) = E((1+a)\theta_{MVU}) = (1+a)E(\theta_{MVU}) = (1+a)\theta_0$ . For  $a \neq 0$ ,  $\theta_b$  is biased estimator of  $\theta_0$ .
- $MSE(\theta_{MVU}) = E[(\theta_{MVU} - E[\theta_{MVU}])^2]$ . The mean square error cannot be zero unless  $\theta_{MVU}$  has zero variance. This cannot be achieved in practice since we have a finite data set of points  $N$ . Zero variance is achieved for infinite size of data samples.
- $MSE(\theta_b) = E\left[\left((1+a)\theta_{MVU} - E[(1+a)\theta_{MVU}]\right)^2\right] + (E[(1+a)\theta_{MVU}] - \theta_0)^2 =$   
 $(1+a)^2 E[\theta_{MVU} - E[\theta_{MVU}]]^2 + (E[\theta_{MVU} - \theta_0 + aE[\theta_{MVU}]]^2 =$   
 $MSE(\theta_{MVU}) + (E[\theta_{MVU}])^2 \Rightarrow$   
 $MSE(\theta_b) = (1+a)^2 MSE(\theta_{MVU}) + (a\theta_0)^2$

e)

$$\begin{aligned} \text{MSE}(\theta_b) < \text{MSE}(\theta_{MVU}) &\Rightarrow (1+a)^2 \text{MSE}(\theta_{MVU}) + a^2 \theta_0^2 < \text{MSE}(\theta_{MVU}) \\ &\Rightarrow (a^2 + 2a) \text{MSE}(\theta_0^2) + 2a \text{MSE}(\theta_{MVU}) < 0 \\ &\Rightarrow a^2 (\text{MSE}(\theta_{MVU}) + \theta_0^2) + 2a \text{MSE}(\theta_{MVU}) < 0 \end{aligned}$$

The above polynomial is negative when  $-\frac{2\text{MSE}(\theta_{MVU})}{\text{MSE}(\theta_{MVU}) + \theta_0^2} < a < 0$ .

f) If we add 1 to the last equality, we have:

$$1 - \frac{2\text{MSE}(\theta_{MVU})}{\text{MSE}(\theta_{MVU}) + \theta_0^2} < a + 1 < 1 \Rightarrow |a + 1| < 1$$

We conclude that:

$$|\theta_b| < |1 + a| |\theta_{MVU}|$$

g)

$$\frac{\partial \text{MSE}(\theta_b)}{\partial a} = 0 \Rightarrow (1+a) \text{MSE}(\theta_{MVU}) + a \theta_0^2 = 0 \Rightarrow a^* = -\frac{\text{MSE}(\theta_{MVU})}{\text{MSE}(\theta_{MVU}) + \theta_0^2}$$

h)  $a^*$  cannot be determined because  $\theta_0$  is unknown.

### Exercise 3

a) From the equation of the LS we have that:

$$N\theta = \sum_{n=1}^N y_n \Rightarrow \theta = \frac{\sum_{n=1}^N y_n}{N}$$

b)  $E[y_n] = E[\theta_0] + E[\eta_n] = \theta_0 + 0 = \theta_0$ ,  $E[\eta_n] = 0$  since the noise is zero mean.c)  $E[\bar{y}] = E\left[\frac{\sum_{n=1}^N y_n}{N}\right] = \frac{N\theta_0}{N} = \theta_0$ .

d)

e) Using the formula from 1a, we have that:

$$(N + \lambda)\hat{\theta} = \sum_{n=1}^N y_n \Rightarrow \hat{\theta} = \frac{\sum_{n=1}^N y_n}{N + \lambda}$$

f) From (a) we have that the ridge regression estimator  $\hat{\theta}$  is related to the LS estimator denoted as  $\theta'_{MVU}$  as follows:

$$\theta = \frac{N}{N + \lambda} \theta'_{MVU}$$

g)

$$E[\hat{\theta}] = E\left[\frac{N}{N+\lambda}\theta'_{MVU}\right] = \frac{N}{N+\lambda}E[\theta'_{MVU}] = \frac{N}{N+\lambda}\theta_0$$

Thus,  $\hat{\theta}$  is a biased estimator.

h) From (f) we have that  $|\theta| = \left|\frac{N}{N+\lambda}\right||\theta'_{MVU}|$  and  $\left|\frac{N}{N+\lambda}\right| < 1 \text{ for } \lambda > 0 \Rightarrow |\hat{\theta}| < |\theta'_{MVU}|$

i) Recalling from exercise 2 that  $\hat{\theta}_b = (1+a)\theta'_{MVU}$  we can get that  $(1+a) = \frac{N}{N+\lambda}$

Then from (2e):

$$\frac{\lambda}{N+\lambda} < \frac{2MSE(\theta'_{MVU})}{MSE(\theta'_{MVU}) + \theta_0^2} \Rightarrow \lambda < \frac{2MSE(\theta'_{MVU})N}{\theta_0^2 - MSE(\theta'_{MVU})}$$

## Exercise 4

d) It is noticed the 8<sup>th</sup> degree polynomial model has high complexity for this problem, which is linear, thus our model overfit the data. It is worth noticing that the coefficients of higher orders are estimated to large values. On the contrary, in the regularized LS the higher orders coefficient takes values close to zero. Regularization is a useful tool when the model which generate the data is unknown.