

## I) Polynomial Regression

$$a) \Phi = \begin{bmatrix} 1 & 0,8 & 0,64 & 0,512 \\ 1 & 1 & 1 & 1 \\ 1 & 1,2 & 1,44 & 1,728 \\ 1 & 1,4 & 1,96 & 2,744 \\ 1 & 1,6 & 2,56 & 4,096 \end{bmatrix}$$

b) Polynomial regression

$$E(w) = \sum_{i=1}^m (t^{(i)} - w^T \phi^{(i)})^2 = (T - Xw)^T (T - Xw)$$

$$\frac{\partial E(w)}{\partial w} = 0 \Leftrightarrow \frac{\partial (T - Xw)^T (T - Xw)}{\partial w} = 0 \Leftrightarrow$$

$$\Leftrightarrow \left( \frac{\partial}{\partial w} (T - Xw)^T \right) \cdot (T - Xw) + (T - Xw)^T \left( \frac{\partial}{\partial w} (T - Xw) \right) = 0 \Leftrightarrow$$

$$\Leftrightarrow (-X^T) (T - Xw) + (T - Xw)^T (-X) = 0 \Leftrightarrow$$

$$\Leftrightarrow (-X^T) (T - Xw) + (-X)^T (T - Xw) = 0 \Leftrightarrow$$

$$\Leftrightarrow -2X^T (T - Xw) = 0 \Leftrightarrow X^T (T - Xw) = 0 \Leftrightarrow$$

$$\Leftrightarrow X^T T - X^T X w = 0 \Leftrightarrow X^T X w = X^T T \Leftrightarrow$$

$$\Leftrightarrow w = (X^T X)^{-1} X^T T$$

$$w = \left( \begin{bmatrix} 1 & 0,8 & 0,64 & 0,512 \\ 1 & 1 & 1 & 1 \\ 1 & 1,2 & 1,44 & 1,728 \\ 1 & 1,4 & 1,96 & 2,744 \\ 1 & 1,6 & 2,56 & 4,096 \end{bmatrix}^T \cdot \begin{bmatrix} 1 & 0,8 & 0,64 & 0,512 \\ 1 & 1 & 1 & 1 \\ 1 & 1,2 & 1,44 & 1,728 \\ 1 & 1,4 & 1,96 & 2,744 \\ 1 & 1,6 & 2,56 & 4,096 \end{bmatrix} \right)^{-1} \cdot \begin{bmatrix} 1 & 0,8 & 0,64 & 0,512 \\ 1 & 1 & 1 & 1 \\ 1 & 1,2 & 1,44 & 1,728 \\ 1 & 1,4 & 1,96 & 2,744 \\ 1 & 1,6 & 2,56 & 4,096 \end{bmatrix}^T \cdot \begin{bmatrix} 24 \\ 20 \\ 10 \\ 13 \\ 12 \end{bmatrix} =$$

$$= \begin{bmatrix} 47,942857 \\ -9,7619046 \\ -41,0714288 \\ 20,833333 \end{bmatrix}$$

$$Y(x, w) = 47,942857 - 9,7619046 \cdot x_1 - 41,0714288 \cdot x_2 + 20,833333 \cdot x_3$$

c) Ridge regression ( $\ell_2$  regularization)

$$\log(\lambda/2) = 0 \Leftrightarrow \lambda/2 = e^0 \Leftrightarrow \lambda/2 = 1 \Leftrightarrow \lambda = 2; \quad \|w\|_2^2 = \sum_{j \in w_2} w_j^2$$

$$E(w) = \frac{1}{2} \sum_{k=1}^m (t_k - w^T \phi_k)^2 + \frac{\lambda}{2} \|w\|_2^2 = (T - Xw)^T (T - Xw) + \lambda w^T w =$$

$$= T^T T - w^T X^T T - T^T X w + w^T X^T X w + \lambda w^T w =$$

$$= T^T T - w^T X^T T - w^T X^T T + w^T X^T X w + w^T \lambda I w =$$

$$= T^T T - 2 w^T X^T T + w^T (X^T X + \lambda I) w$$

Now, to search for the  $w$  that minimizes according to Ridge:

$$\frac{\partial E(w)}{\partial w} = -2X^T T + 2(X^T X + \lambda I)w = 0 \Leftrightarrow$$

$$\Leftrightarrow (X^T X + \lambda I)w = X^T T \Rightarrow w = (X^T X + \lambda I)^{-1} X^T T$$

$$w = \left( \begin{bmatrix} 1 & 0,8 & 0,64 & 0,512 \\ 1 & 1 & 1 & 1 \\ 1 & 1,2 & 1,44 & 1,728 \\ 1 & 1,4 & 1,96 & 2,744 \\ 1 & 1,6 & 2,56 & 4,096 \end{bmatrix}^T \cdot \begin{bmatrix} 1 & 0,8 & 0,64 & 0,512 \\ 1 & 1 & 1 & 1 \\ 1 & 1,2 & 1,44 & 1,728 \\ 1 & 1,4 & 1,96 & 2,744 \\ 1 & 1,6 & 2,56 & 4,096 \end{bmatrix} + 2 \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \right)^{-1} \cdot \begin{bmatrix} 1 & 0,8 & 0,64 & 0,512 \\ 1 & 1 & 1 & 1 \\ 1 & 1,2 & 1,44 & 1,728 \\ 1 & 1,4 & 1,96 & 2,744 \\ 1 & 1,6 & 2,56 & 4,096 \end{bmatrix}^T \cdot \begin{bmatrix} 24 \\ 20 \\ 10 \\ 13 \\ 12 \end{bmatrix} =$$

$$= \begin{bmatrix} 7,0450759 \\ 4,64092765 \\ 1,96734046 \\ -1,300881417 \end{bmatrix}$$

$$Y(x, w) = 7,0450759 + 4,64092765 \cdot x_1 + 1,96734046 \cdot x_2 - 1,300881417 \cdot x_3$$

d) LASSO ( $\ell_1$  regularization)

$$E(w) = \frac{1}{2} \sum_{k=1}^m (t_k - w^T \phi_k)^2 + \lambda \|w\|_1 = (T - Xw)^T (T - Xw) + \lambda \|w\|_1 =$$

$$= T^T T - w^T X^T T - T^T X w + w^T X^T X w + \lambda \|w\|_1 =$$

$$= T^T T - w^T X^T T - w^T X^T T + w^T X^T X w + \lambda \|w\|_1 =$$

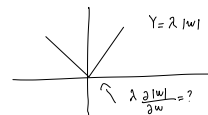
$$= T^T T - 2 w^T X^T T + w^T (X^T X) w + \lambda \|w\|_1$$

$$\frac{\partial E(w)}{\partial w} = -2X^T T + 2(X^T X)w + \left( \frac{\partial \lambda \|w\|_1}{\partial w} \right) = 0$$

→ this is not differentiable  
↙  
So we cannot use differentiation for a closed form of Lasso ( $\ell_1$  regularization)

Lasso ( $\ell_1$  regularization) lacks a closed form solution.

This happens because  $\gamma = \lambda \|w\|_1$  is a non-differentiable function in  $w = 0$



# 11) Neural Network

a)

$$w_1 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 \end{bmatrix} \quad b_1 = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \quad f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$w_2 = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} \quad b_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \eta = 0.1$$

$$w_3 = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad b_3 = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad f(x) = \tanh(x)$$

$$x_0 = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \quad t = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

$$z_1 = w_1 \cdot x_0 + b_1 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 6 \\ 0 \\ 6 \end{bmatrix}$$

$$x_1 = f(z_1) = \begin{bmatrix} f(6) \\ f(0) \\ f(6) \end{bmatrix} = \begin{bmatrix} \tanh(6) \\ \tanh(0) \\ \tanh(6) \end{bmatrix} = \begin{bmatrix} 0.9999 \\ 0 \\ 0.9999 \end{bmatrix}$$

$$z_2 = w_2 \cdot x_1 + b_2 = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} 0.9999 \\ 0 \\ 0.9999 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 3.76157 \\ 3.76157 \end{bmatrix}$$

$$x_2 = f(z_2) = \begin{bmatrix} \tanh(3.76157) \\ \tanh(3.76157) \end{bmatrix} = \begin{bmatrix} 0.99892 \\ 0.99892 \end{bmatrix}$$

$$z_3 = w_3 \cdot x_2 + b_3 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$x_3 = \tanh\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}\right) = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Now, we want to do the backward phase

$$E(t, x_1) = \frac{1}{2} (x_1 - t)^2 \quad \text{Squared error loss}$$

In general we will need to know how to derive all functions in our network.

Let's compute them beforehand:

$$\frac{\partial E}{\partial x_1}(t, x_1) = x_1 - t$$

$$\frac{\partial x_1}{\partial z_1}(z_1) = 1 - \tanh^2(z_1)$$

$$\frac{\partial z_1}{\partial w_1}(w_1, b_1, x_{0,1}) = x_{0,1}$$

$$\frac{\partial z_1}{\partial b_1}(w_1, b_1, x_{0,1}) = 1$$

$$\frac{\partial z_1}{\partial x_{0,1}}(w_1, b_1, x_{0,1}) = w_1$$

Start the recursion, we need the delta from the past layer.

$$\delta_3 = \frac{\partial E}{\partial x_3} \cdot \frac{\partial x_3}{\partial z_3} = (x_3 - t) \cdot (1 - \tanh^2(z_3)) =$$

$$= ([0] - [-1]) \cdot ([1] - [0]) = [1]$$

We can use the recursion to compute the delta from the hidden layers.

$$\delta_2 = \frac{\partial z_2}{\partial x_2} \cdot \delta_3 \cdot \frac{\partial x_2}{\partial z_2} = (w_3)^T \cdot \delta_3 \cdot (1 - \tanh^2(z_2)) =$$

$$= \begin{bmatrix} 0 & 0 \end{bmatrix} \cdot [1] \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 0.99892^2 \\ 0.99892^2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\delta_1 = \frac{\partial z_1}{\partial x_1} \cdot \delta_2 \cdot \frac{\partial x_1}{\partial z_1} = (w_2)^T \cdot \delta_2 \cdot (1 - \tanh^2(z_1)) =$$

$$= \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 0 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 0.9999^2 \\ 0.9999^2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

Finally, we can go to the last phase and perform the updates. We start with the first layer.

$$\frac{\partial E}{\partial w_1} = \delta_1 \cdot \frac{\partial z_1}{\partial w_1} = \delta_1 \cdot (x_0)^T = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \cdot [1, 1, 1, 1] = \begin{bmatrix} 0 & 0 & 0 & 0 \end{bmatrix}$$

$$w_1 = w_1 - \eta \frac{\partial E}{\partial w_1} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 \end{bmatrix} - 0.1 \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 \end{bmatrix}$$

$$\frac{\partial E}{\partial b_1} = \delta_1 \cdot \frac{\partial z_1}{\partial b_1} = \delta_1 \cdot [1] = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$b_1 = b_1 - \eta \frac{\partial E}{\partial b_1} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} - 0.1 \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

Now the second:

$$\frac{\partial E}{\partial w_2} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$$

$$w_2 = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

$$\frac{\partial E}{\partial b_2} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$b_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\frac{\partial E}{\partial w_3} = \begin{bmatrix} -1 \\ 1 \end{bmatrix} \cdot [0.99892 \quad 0.99892] =$$

$$= \begin{bmatrix} -0.99892 & 0.99892 \\ 0.99892 & -0.99892 \end{bmatrix}$$

$$w_3 = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} - 0.1 \cdot \begin{bmatrix} -0.99892 & 0.99892 \\ 0.99892 & -0.99892 \end{bmatrix} =$$

$$= \begin{bmatrix} 0.099892 & -0.099892 \\ -0.099892 & 0.099892 \end{bmatrix}$$

$$\frac{\partial E}{\partial b_3} = \delta_2 = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad b_3 = b_3 - 0.1 \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

b) We cannot use Cross Entropy Error in the previous exercise because the  $t \in [0, 1]$ . The formula gives us a probability distribution and probabilities  $\in [0, 1]$ . In the previous exercise our target is  $t = (1 - t)^T$  and  $-1 \notin [0, 1]$ , so we cannot apply Cross Entropy Error.

c)  $t = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad x_0 = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$

Because we only apply Softmax in the output unit, we can reuse the computed results from exercise a) as follows:

$$x_0 = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \quad z_1 = \begin{bmatrix} 6 \\ 0 \\ 6 \end{bmatrix} \quad x_1 = \begin{bmatrix} 0.9999 \\ 0.76157 \\ 0.9999 \end{bmatrix} \quad z_2 = \begin{bmatrix} 3.76157 \\ 3.76157 \end{bmatrix}$$

$$x_2 = \begin{bmatrix} 0.99892 \\ 0.99892 \end{bmatrix} \quad z_3 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

we only use softmax on the  $x_3$ :  $x_3 = \text{softmax}([0]) = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$

Because we applied Softmax on the output unit, we now can use Cross Entropy Loss because we no longer have 0 on  $x_3 = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$

$$E(t, x_3) = - \sum_{i=1}^2 t_i \log x_{3,i}$$

$$\delta_{3,i} = x_{3,i} - t_i$$

The remaining derivatives can be computed as before from exercise

a):

$$\frac{\partial x_1}{\partial z_1}(z_1) = 1 - \tanh^2(z_1)$$

$$\frac{\partial z_1}{\partial w_1}(w_1, b_1, x_{0,1}) = x_{0,1}$$

$$\frac{\partial z_1}{\partial b_1}(w_1, b_1, x_{0,1}) = 1$$

$$\frac{\partial z_1}{\partial w_1}(w_1, b_1, x_{0,1}) = 1$$

$$\frac{\partial z_1}{\partial b_1}(w_1, b_1, x_{0,1}) = 1$$

To start the recursion, we need the delta from the last layer

$$\delta_3 = (x_3 - t) = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} - \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} -0.5 \\ 0.5 \end{bmatrix}$$

$$\delta_2 = \frac{\partial z_2}{\partial x_2} \cdot \delta_3 \cdot \frac{\partial x_2}{\partial z_2} = (w_3)^T \cdot \delta_3 \cdot (1 - \tanh^2(z_2)) = \begin{bmatrix} 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} -0.5 \\ 0.5 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 0.99892^2 \\ 0.99892^2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\delta_1 = \frac{\partial z_1}{\partial x_1} \cdot \delta_2 \cdot \frac{\partial x_1}{\partial z_1} = (w_2)^T \cdot \delta_2 \cdot (1 - \tanh^2(z_1)) = \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 0 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 0.9999^2 \\ 0.9999^2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

Finally, we can perform the updates (the last phase). We will start with the first layer:

$$\frac{\partial E}{\partial w_1} = \delta_1 \cdot \frac{\partial z_1}{\partial w_1} = \delta_1 \cdot (x_0)^T = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \cdot [1, 1, 1, 1] = \begin{bmatrix} 0 & 0 & 0 & 0 \end{bmatrix}$$

$$w_1 = w_1 - \eta \frac{\partial E}{\partial w_1} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 \end{bmatrix} - 0.1 \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 \end{bmatrix}$$

$$\frac{\partial E}{\partial b_1} = \delta_1 \cdot \frac{\partial z_1}{\partial b_1} = \delta_1 \cdot [1] = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad b_1 = b_1 - \eta \frac{\partial E}{\partial b_1} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} - 0.1 \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

Now the second:

$$\frac{\partial E}{\partial w_2} = \delta_2 \cdot \frac{\partial z_2}{\partial w_2} = \delta_2 \cdot (x_1)^T = \begin{bmatrix} 0 & 0 \end{bmatrix} \cdot [0.9999 \quad 0.76157 \quad 0.9999] = \begin{bmatrix} 0 & 0 \end{bmatrix}$$

$$w_2 = w_2 - \eta \frac{\partial E}{\partial w_2} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} - 0.1 \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

$$\frac{\partial E}{\partial b_2} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad b_2 = b_2 - \eta \frac{\partial E}{\partial b_2} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

All that is left is to update the parameters for the output layer:

$$\frac{\partial E}{\partial w_3} = \delta_3 \cdot \frac{\partial z_3}{\partial w_3} = \delta_3 \cdot (x_2)^T = \begin{bmatrix} -0.5 \\ 0.5 \end{bmatrix} \cdot [0.99892 \quad 0.99892] = \begin{bmatrix} -0.49946 & 0.49946 \end{bmatrix}$$

$$w_3 = w_3 - \eta \frac{\partial E}{\partial w_3} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} - 0.1 \begin{bmatrix} -0.49946 & 0.49946 \\ 0.49946 & -0.49946 \end{bmatrix} = \begin{bmatrix} 0.049946 & -0.049946 \\ -0.049946 & 0.049946 \end{bmatrix}$$

$$\frac{\partial E}{\partial b_3} = \delta_3 = \begin{bmatrix} -0.5 \\ 0.5 \end{bmatrix}$$

$$b_3 = b_3 - \eta \frac{\partial E}{\partial b_3} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} - 0.1 \begin{bmatrix} -0.5 \\ 0.5 \end{bmatrix} = \begin{bmatrix} 0.05 \\ -0.05 \end{bmatrix}$$