

# Homework III

Madalena Yang 110206  
Madalena Korta 110355

1.	D	$y_1$	$y_2$	$y_{sum}$	$\Phi(y_1, y_2) = y_1 \times y_2$
	$x_1$	2	2	3,5	4
	$x_2$	1	1	1,0	1
	$x_3$	3	2	3,8	6
	$x_4$	6	3	10,1	18
	$x_5$	8	1	8,5	8

OLS closed form solution  
 $W = (X^T X)^{-1} X^T z$

$$X = \begin{bmatrix} 1 & 4 \\ 1 & 1 \\ 1 & 6 \\ 1 & 18 \\ 1 & 8 \end{bmatrix}$$

$$z = \begin{bmatrix} 3,5 \\ 1,0 \\ 3,8 \\ 10,1 \\ 8,5 \end{bmatrix}$$

$$W = [1,46136 \quad 0,52955]^T$$

(calculated in code)

2. This time, we'll learn a Ridge Regression, so the weights are given by  $W = (X^T X + \lambda I)^{-1} X^T z$ , with  $\lambda = 1$  penalty factor

Ridge  
 $W = (X^T X + \lambda I)^{-1} X^T z$

$$W = [0,97319 \quad 0,56921]^T \rightarrow \text{calculated in code}$$

A flexible model is very susceptible to noise, which means it fits the training data well but fails to generalize. Regularization helps prevent overfitting by restricting the model, penalizing large and unstable coefficients. This results in more generalized, well-fitted and robust models.

As expected, after adding Ridge Regularization, the slope barely changed and the intercept decreased from 1,46136 to 0,97319. This is because Ridge penalizes large coefficients to reduce overfitting and to prevent them from influencing predictions too heavily compared to other more important ones.

In conclusion, the results obtained with Regularization were expected = larger coefficients shrunk so they don't influence the predictions as heavily, and the importance of each was calibrated.

3.

OLS model:

$$W = [1,46136 \quad 0,52955]^T$$

$$\hat{z}_i = w_0 + w_1 x_i$$

$$\hat{z}_i = 1,46136 + 0,52955 x_i$$

$$\hat{z} = f(x) = w_0 + w_1 x_1 + \dots + w_n x_n$$

Mean absolute error

$$MAE = \frac{1}{n} \sum_{i=1}^n |\hat{z}_i - z_i|$$

↑  
prediction  
↓  
 $y_{sum}$



D	$y_1$	$y_2$	$y_{num}$	$y_1 \times y_2$	$\hat{z}_i = w_0 + w_1 x_i$	
$u_1$	2	2	3,5	4	3,57956	TRAIN
$u_2$	1	1	1,0	1	1,99091	
$u_3$	3	2	3,8	6	4,63866	
$u_4$	6	3	10,1	18	10,99326	
$u_5$	8	1	8,5	8	5,69776	
$u_6$	0	2	1,0	0	1,46136	TEST
$u_7$	3	4	6,2	12	7,81596	
$u_8$	5	1	3,6	5	4,10911	

• Train MAE:  $n=5$  samples

$$z = [3,5 \quad 1,0 \quad 3,8 \quad 10,1 \quad 8,5]^T$$

$$MAE_{\text{out train}} = \frac{1}{5} (|3,57956 - 3,5| + |1,99091 - 1,0| + |4,63866 - 3,8| + |10,99326 - 10,1| + |5,69776 - 8,5|) \approx 1,12093$$

• Test MAE:  $n=3$  samples

$$z = [1,0 \quad 6,2 \quad 3,6]^T$$

$$MAE_{\text{out test}} = \frac{1}{3} (|1,46136 - 1,0| + |7,81596 - 6,2| + |4,10911 - 3,6|) \approx 0,86214$$

Ridge model:  $w = [0,97319 \quad 0,56921]^T$   $\hat{z}_i = 0,97319 + 0,56921 x_i$

D	$y_1 \times y_2$	$\hat{z}_i$
$u_1$	4	3,25003
$u_2$	1	1,54240
$u_3$	6	4,38845
$u_4$	18	11,21897
$u_5$	8	5,52687
$u_6$	0	0,97319
$u_7$	12	7,80371
$u_8$	5	3,81924



• Train MAE :  $n=5$  samples

$$z = [3.5 \quad 1.0 \quad 3.8 \quad 10.1 \quad 8.5]^T$$

$$MAE_{\text{Ridge train}} = \frac{1}{5} (|3.25003 - 3.5| + |1.54240 - 1.0| + |4.38845 - 3.8| + |11.21897 - 10.1| + |5.52687 - 8.5|) \\ \approx 1.09458$$

• Test MAE :  $n=3$  samples

$$z = [1.0 \quad 6.2 \quad 3.6]^T$$

$$MAE_{\text{Ridge test}} = \frac{1}{3} (|10.97319 - 1.0| + |17.80371 - 6.2| + |13.81924 - 3.6|) \approx 0.61659$$

As expected, both training and testing MAE decreased for the Ridge Regression model compared to OLS. In Ridge, the model achieves lower variance, which means it is more stable and less sensitive to noise, at the cost of a little bit more bias. This bias variance tradeoff improves model robustness and stability.

We observed a better generalization of data with Ridge, and even a slightly better train fit. This means Ridge regularized coefficients and removed some instability.

4.

$$W^{[1]} = \begin{bmatrix} 0.1 & 0.1 \\ 0.1 & 0.2 \\ 0.2 & 0.1 \end{bmatrix}$$

$$b^{[1]} = \begin{bmatrix} 0.1 \\ 0 \\ 0.1 \end{bmatrix}$$

$$W^{[2]} = \begin{bmatrix} 1 & 2 & 2 \\ 1 & 2 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

$$b^{[2]} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

$$x_1 = [2 \quad 2]^T$$

$$\begin{array}{ccc} & 0 & 0 \\ & 0 & 0 \\ & 0 & 0 \\ & 0 & 0 \\ L_0 & 1 & 2 \\ & \searrow & \searrow \\ & \text{fe=softmax} & \text{fx=sigmoid} \end{array}$$

1. Forward Propagation

$$z^{[1]} = W^{[1]} x^{[0]} + b^{[1]}$$

$$k^{[1]} = \text{fa}(z^{[1]})$$

activation function

1. Forward Propagation :  $z^{[1]} \rightarrow k^{[1]} \rightarrow z^{[2]} \rightarrow k^{[2]}$   
output

$$z^{[1]} = W^{[1]} x^{[0]} + b^{[1]} = \begin{bmatrix} 0.1 & 0.1 \\ 0.1 & 0.2 \\ 0.2 & 0.1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix} + \begin{bmatrix} 0.1 \\ 0 \\ 0.1 \end{bmatrix} = \begin{bmatrix} 0.5 \\ 0.6 \\ 0.7 \end{bmatrix}$$

$$k^{[1]} = \text{fa}(z^{[1]}) = \text{sigmoid}(z^{[1]}) = \sigma(z^{[1]}) = \frac{1}{1 + e^{-z^{[1]}}} = \begin{bmatrix} 0.62246 \\ 0.64566 \\ 0.66819 \end{bmatrix}$$



$$z^{[2]} = w^{[2]} k^{[1]} + b^{[2]} = \begin{bmatrix} 1 & 2 & 2 \\ 1 & 2 & 1 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 0,62246 \\ 0,64566 \\ 0,66819 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 4,25016 \\ 3,58197 \\ 2,93631 \end{bmatrix}$$

$$u^{[2]} = \text{fa}(z^{[2]}) = \text{softmax}(z^{[2]}) = \text{softmax} \left( \begin{bmatrix} 4,25016 \\ 3,58197 \\ 2,93631 \end{bmatrix} \right) \approx \begin{bmatrix} 0,56135 \\ 0,28777 \\ 0,15088 \end{bmatrix}$$

2. Back Propagation:  $\delta^{[2]} \rightarrow \delta^{[1]}$  Stochastic: one observation

Cross-entropy loss:  $-\sum_{i=1}^N \sum_{c=1}^{|C|} t_c^{(i)} \log(u^{(i)})$  where  $t$  represents the target classes  $\{A, B, C\}$

in this case because we're in the last layer, so AF is softmax

$$\delta^{[2]} = \frac{\partial \mathcal{E}}{\partial u^{[2]}} \frac{\partial u^{[2]}}{\partial z^{[2]}} = u^{[2]} - t^{[2]} =$$

$$= \begin{bmatrix} 0,56135 \\ 0,28777 \\ 0,15088 \end{bmatrix} - \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} -0,43865 \\ 0,28777 \\ 0,15088 \end{bmatrix}$$

the target class is A

$$\delta^{[1]} = (w^{[2]})^T \delta^{[2]} \left( \frac{\partial u^{[2]}}{\partial z^{[2]}} \right) = (w^{[2]})^T \delta^{[2]} \cdot u^{[1]} \cdot (1 - u^{[1]}) =$$

$$= \begin{bmatrix} 1 & 2 & 2 \\ 1 & 2 & 1 \\ 1 & 1 & 1 \end{bmatrix}^T \begin{bmatrix} -0,43865 \\ 0,28777 \\ 0,15088 \end{bmatrix} \cdot \begin{bmatrix} 0,62246 \\ 0,64566 \\ 0,66819 \end{bmatrix} \cdot \begin{bmatrix} 0,37754 \\ 0,35434 \\ 0,33181 \end{bmatrix} = \begin{bmatrix} 0 \\ -0,03452 \\ -0,09725 \end{bmatrix}$$

3. Update weights ( $\eta = 0,5$ )

$$w^{[2]} = w^{[2]} - \eta \frac{\partial \mathcal{E}}{\partial w^{[2]}} = w^{[2]} - \eta \delta^{[2]} (k^{[1]})^T =$$

$$= \begin{bmatrix} 1 & 2 & 2 \\ 1 & 2 & 1 \\ 1 & 1 & 1 \end{bmatrix} - 0,5 \begin{bmatrix} -0,43865 \\ 0,28777 \\ 0,15088 \end{bmatrix} \begin{bmatrix} 0,62246 \\ 0,64566 \\ 0,66819 \end{bmatrix}^T =$$

$$= \begin{bmatrix} 1,13652 & 2,14161 & 2,14655 \\ 0,91044 & 1,90710 & 0,90386 \\ 0,95304 & 0,95129 & 0,94959 \end{bmatrix}$$

2. Back Propagation

$$\delta^{[1]} = \frac{\partial \mathcal{E}}{\partial z^{[1]}}$$

$$\text{last layer: } \delta^{[1]} = \frac{\partial \mathcal{E}}{\partial u^{[1]}} \frac{\partial u^{[1]}}{\partial z^{[1]}}$$

$$\text{other: } \delta^{[1]} = \left( \frac{\partial z^{[2]}}{\partial u^{[1]}} \right)^T \delta^{[2]} \left( \frac{\partial u^{[1]}}{\partial z^{[1]}} \right)$$

derivative of AF

$$\delta^{[1]} = \frac{\partial \mathcal{E}}{\partial z^{[1]}} = \frac{\partial \mathcal{E}}{\partial u^{[1]}} \frac{\partial u^{[1]}}{\partial z^{[1]}}$$

Sigmoid: (for hidden layers in this case)

$$\frac{\partial u^{[1]}(z^{[1]})}{\partial z^{[1]}} = \sigma(z^{[1]}) (1 - \sigma(z^{[1]})) =$$

$$= u^{[1]} (1 - u^{[1]})$$

3. Update weights

$$w^{[1]} = w^{[1]} - \eta \frac{\partial \mathcal{E}}{\partial w^{[1]}}$$

$$(b^{[1]}) = b^{[1]} - \eta \frac{\partial \mathcal{E}}{\partial b^{[1]}}$$

$$\frac{\partial \mathcal{E}}{\partial w^{[1]}} = \frac{\partial \mathcal{E}}{\partial z^{[1]}} \frac{\partial z^{[1]}}{\partial w^{[1]}} =$$

$$= \delta^{[1]} \frac{\partial z^{[1]}}{\partial w^{[1]}} = \delta^{[1]} x^{[0]}$$

$$\frac{\partial \mathcal{E}}{\partial b^{[1]}} = \frac{\partial \mathcal{E}}{\partial z^{[1]}} \frac{\partial z^{[1]}}{\partial b^{[1]}} = \delta^{[1]}$$



$$b^{[2]} = b^{[1]} - \eta \frac{\partial \mathcal{L}}{\partial b^{[1]}} = b^{[1]} - \eta f^{[1]} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} - 0,5 \begin{bmatrix} -0,43865 \\ 0,28777 \\ 0,15088 \end{bmatrix} = \begin{bmatrix} 1,21933 \\ 0,85612 \\ 0,92456 \end{bmatrix}$$

$$w^{[1]} = w^{[1]} - \eta \frac{\partial \mathcal{L}}{\partial w^{[1]}} = w^{[1]} - \eta \delta^{[1]} (x^{[0]})^T = \quad \quad \quad x^{[0]} = \begin{bmatrix} 2 & 2 \end{bmatrix}$$

$$= \begin{bmatrix} 0,1 & 0,1 \\ 0,1 & 0,2 \\ 0,2 & 0,1 \end{bmatrix} - 0,5 \begin{bmatrix} 0 \\ -0,03452 \\ -0,09725 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix}^T =$$

$$= \begin{bmatrix} 0,1 & 0,1 \\ 0,13452 & 0,23452 \\ 0,29725 & 0,19725 \end{bmatrix}$$

$$b^{[1]} = b^{[1]} - \eta \frac{\partial \mathcal{L}}{\partial b^{[1]}} = b^{[1]} - \eta f^{[1]} = \begin{bmatrix} 0,1 \\ 0 \\ 0,1 \end{bmatrix} - 0,5 \begin{bmatrix} 0 \\ -0,03452 \\ -0,09725 \end{bmatrix} = \begin{bmatrix} 0,1 \\ 0,01726 \\ 0,14863 \end{bmatrix}$$

### Explanation:

A model with no activation function simply outputs a linear combination of its inputs. Since each layer only performs a linear transformation, as a result, the overall model can only represent linear relationships and fails to capture more complex, nonlinear patterns in the data.

By contrast, a sigmoid activation function introduces nonlinearity, enabling the neural network to learn more complex functions and nonlinear decision boundaries.

Additionally, the sigmoid function squeezes inputs into the (0,1) range, effectively regularizing the output. This allows the model to interpret the output as a probability, which makes classification more intuitive and the results easier to interpret. The compression effect of the sigmoid also smooths the influence of extreme inputs, making the model's output more stable and less sensitive to very large values or even outliers.