

HW - 3

[I]

1.

$$\text{OLS} : \mathbf{w}_{\text{OLS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\mathbf{y} = \begin{bmatrix} 3.5 \\ 1.0 \\ 3.8 \\ 10.1 \\ 8.5 \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 0.5 & 0.3 & 0.4 & 0.5 \\ 1 & 0.5 & 0.3 & 0.4 & 0.5 \\ 1 & 1 & 6 & 18 & 8 \end{bmatrix} = \begin{bmatrix} 1 & 4 \\ 1 & 1 \\ 1 & 6 \\ 1 & 18 \\ 1 & 8 \end{bmatrix}$$

$$\mathbf{X}^T \mathbf{X} = \begin{bmatrix} 5 & 37 \\ 37 & 441 \end{bmatrix} \quad \mathbf{X}^T \mathbf{y} = \begin{bmatrix} 26.9 \\ 287.6 \end{bmatrix}$$

$$\det(\mathbf{X}^T \mathbf{X}) = 836 \quad (\mathbf{X}^T \mathbf{X})^{-1} = \frac{1}{836} \begin{bmatrix} 441 & -37 \\ -37 & 5 \end{bmatrix}$$

$$\mathbf{w}_{\text{OLS}} = \frac{1}{836} \begin{bmatrix} 441 & -37 \\ -37 & 5 \end{bmatrix} \begin{bmatrix} 26.9 \\ 287.6 \end{bmatrix}$$

$$= \begin{bmatrix} 1,46136 \\ 0,52954 \end{bmatrix}$$

$$\hat{\mathbf{y}}_{\text{OLS}} = 1,46136 + 0,52954(y_1, y_2)$$

$$2. \quad \Lambda = \begin{bmatrix} 0 & 0 \\ 0 & X \end{bmatrix} - \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \quad \det[(X^T X + \Lambda)^{-1}] = 841$$

$$\text{Ridge: } w_{\text{ridge}} = (X^T X + \Lambda)^{-1} X^T y$$

$$w_{\text{ridge}} = (X^T X + \Lambda)^{-1} X^T y = \frac{1}{841} \begin{bmatrix} 441 & -37 \\ -37 & 5 \end{bmatrix} \begin{bmatrix} 26.7 \\ 287.6 \end{bmatrix}$$

$$= \begin{bmatrix} 1.48466 \\ 0.526391 \end{bmatrix}$$

$$\hat{y}_{\text{Ridge}} = 1.48466 + 0.526391(y_1, y_2)$$

Comparison: With Ridge the Coefficients shrink due to regularization (large Coefficients get pulled towards zero).

3.

Train

$$\bullet \text{OLS: } \hat{y}_{\text{OLS}} = 1,46136 + 0,52954(y_1, y_2)$$

$$\bullet \text{Ridge: } \hat{y}_{\text{ridge}} = 1,48466 + 0,52659(y_1, y_2)$$

$$\text{Error} = y_{\text{true}} - \hat{y}_{\text{predicted}}$$

| x | y, y_1 | \hat{y}_{OLS} | \hat{y}_{Ridge} | $ \text{Error} _{\text{OLS}}$ | $ \text{Error} _{\text{Ridge}}$ |
|-------|----------|------------------------|--------------------------|-------------------------------|---------------------------------|
| x_1 | 4 | 3,57952 | 3,598248 | 0,07952 | 0,09025 |
| x_2 | 7 | 1,99090 | 2,011057 | 0,99090 | 1,011057 |
| x_3 | 6 | 4,63860 | 4,643042 | 0,83860 | 0,84304 |
| x_4 | 18 | 10,99308 | 10,95981 | 0,89308 | 0,85981 |
| x_5 | 8 | 5,69768 | 5,69584 | 2,80232 | 2,80416 |

$$\text{MAE}_{\text{train, OLS}} = \frac{0,07952 + 0,99090 + 0,83860 + 0,89308 + 2,80232}{5} \approx 1,12088$$

$$\text{MAE}_{\text{train, ridge}} = \frac{0,09025 + 1,011057 + 0,84304 + 0,85981 + 2,80416}{5} \approx 1,12066$$

Test

| x | y, y_1 | \hat{y}_{OLS} | \hat{y}_{Ridge} | $ \text{Error} _{\text{OLS}}$ | $ \text{Error} _{\text{Ridge}}$ |
|-------|----------|------------------------|--------------------------|-------------------------------|---------------------------------|
| x_6 | 0 | 1,46136 | 1,48466 | 0,46136 | 0,48466 |
| x_7 | 12 | 7,81584 | 7,80142 | 1,61584 | 1,601424 |
| x_8 | 5 | 4,10906 | 4,11665 | 0,50906 | 0,51665 |

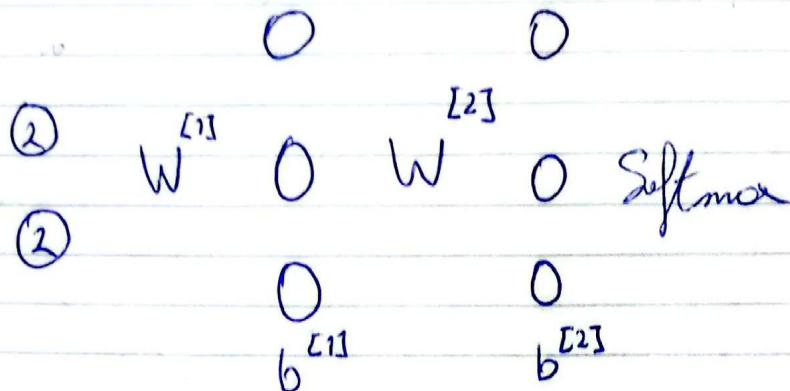
$$\text{MAE}_{\text{test, OLS}} = \frac{0,46136 + 1,61584 + 0,50906}{3} \approx 0,86207$$

$$\text{MAE}_{\text{test, ridge}} = \frac{0,48466 + 1,601424 + 0,51665}{3} \approx 0,86758$$

Explanation: The Ridge model shows a slightly higher training MAE and a similar test MAE compared to OLS. This behavior is what we expect from regularization (Ridge adds a small bias aiming to reduce overfitting and improve stability on new data [although, that is not prominent with these results]).

HW - III

4.



Forward propagation:

$$Z^{[i]} = w^{[i]} \times X^{[i-1]} + b^{[i]}$$

1st layer: $Z^{[1]} = \begin{bmatrix} 0,1 & 0,1 \\ 0,1 & 0,2 \\ 0,2 & 0,1 \end{bmatrix} \cdot \begin{bmatrix} 2 \\ 2 \end{bmatrix} + \begin{bmatrix} 0,1 \\ 0 \\ 0,1 \end{bmatrix} = \begin{bmatrix} 0,4 \\ 0,6 \\ 0,6 \end{bmatrix} + \begin{bmatrix} 0,1 \\ 0 \\ 0,1 \end{bmatrix}$

$$= \begin{bmatrix} 0,5 \\ 0,6 \\ 0,7 \end{bmatrix}$$

Sigmoid: $\delta(n) = \frac{e^n}{1 + e^n}$

Output: $\delta \left(\begin{bmatrix} 0,5 \\ 0,6 \\ 0,7 \end{bmatrix} \right) = \begin{bmatrix} 0,62246 \\ 0,64566 \\ 0,66819 \end{bmatrix}$

$$Z^{[2]} = \begin{bmatrix} 1 & 2 & 2 \\ 1 & 2 & 1 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 0,62246 \\ 0,64566 \\ 0,66819 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 3,2502 \\ 2,58197 \\ 1,93631 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

$$= \begin{bmatrix} 4,2502 \\ 3,58197 \\ 2,93631 \end{bmatrix} \rightarrow \text{Softmax}$$

$$\text{softmax: } \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

$$e^{-4.2502} + e^{-3.58197} + e^{-2.93631} \approx 124.9099$$

$$\text{Find } \mathbf{x}^{[2]} : \begin{bmatrix} 0.56136 \\ 0.28776 \\ 0.150878 \end{bmatrix}$$

Back propagation:

$$\text{Onde } x_1 \rightarrow t = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \begin{matrix} A \\ B \\ C \end{matrix}$$

Actualización: $\delta^{[1]} \quad \delta^{[2]}$

$$W_{\text{new}} = W - \eta \frac{\partial E}{\partial W^{[2]}}$$

$$\frac{\partial E}{\partial W^{[2]}} = \delta^{[1]} (x^{[1]})^T \quad \frac{\partial E}{\partial b^{[2]}} = \delta^{[2]}$$

Output layer:

$$\delta^{[2]} = g' - t = \delta^{[1]} \begin{pmatrix} 4.2502 \\ 3.58197 \\ 2.93631 \end{pmatrix} - \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0.56136 \\ 0.28776 \\ 0.150878 \end{bmatrix} - \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} -0.43864 \\ 0.28776 \\ 0.150878 \end{bmatrix} = \frac{\partial E}{\partial b^{[2]}}$$

$$\frac{\partial E}{\partial W^{[2]}} = \begin{bmatrix} -0.43864 \\ 0.28776 \\ 0.150878 \end{bmatrix} \cdot \begin{bmatrix} 0.62246 & 0.64566 & 0.66819 \end{bmatrix}$$

$$= \begin{bmatrix} -0.27303 & -0.283212 & -0.293045 \\ 0.17912 & 0.185795 & 0.192278 \\ 0.093915 & 0.097416 & 0.1008145 \end{bmatrix}$$

[2]

$$W_{\text{new}} = \begin{bmatrix} 1 & 2 & 2 \\ 1 & 2 & 1 \\ 1 & 1 & 1 \end{bmatrix} - 0.5 \begin{bmatrix} -0.21303 & 0.283212 & -0.293093 \\ 0.17912 & 0.185795 & 0.192278 \\ 0.093915 & 0.097416 & 0.100815 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 2 & 2 \\ 1 & 2 & 1 \\ 1 & 1 & 1 \end{bmatrix} - \begin{bmatrix} -0.136515 & -0.14161 & 0.14655 \\ 0.68956 & 0.092898 & 0.09614 \\ 0.046958 & 0.048708 & 0.05041 \end{bmatrix}$$

$$\approx \begin{bmatrix} 1.136515 & 2.14161 & 2.14655 \\ 0.91044 & 1.907102 & 0.90386 \\ 0.953042 & 0.951292 & 0.94959 \end{bmatrix}$$

[2]

[2]

$$b_{\text{new}} = b - \gamma \frac{\partial E}{\partial b^{[2]}} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} - 0.5 \begin{bmatrix} -0.43864 \\ 0.28776 \\ 0.150878 \end{bmatrix}$$

$$= \begin{bmatrix} 1.21932 \\ 0.85612 \\ 0.924561 \end{bmatrix}$$

[1]

[1]

$$W_{\text{new}} = W - \gamma \frac{\partial E}{\partial W^{[1]}} A, \quad \frac{\partial E}{\partial W^{[1]}} = \delta^{[1]} (x^{[0]})^T$$

In this layer:

$$\delta^{[1]} = (W_{\text{new}}^{[2]})^T \delta^{[2]} \cdot f'(z^{[1]}) \quad f' \rightarrow \text{derivative of sigmoid}$$

$$f'(z) = g(z)(1 - g(z))$$

$$\delta^{[1]} = \begin{bmatrix} 1.136515 & 2.14161 & 2.14655 \\ 0.91044 & 1.907102 & 0.90386 \\ 0.953042 & 0.951292 & 0.94959 \end{bmatrix}^T \begin{bmatrix} -0.43864 \\ 0.28776 \\ 0.150878 \end{bmatrix} \cdot f'(z^{[1]})$$

$$= \begin{bmatrix} -0.09274 \\ -0.24708 \\ -0.53819 \end{bmatrix} \cdot \begin{bmatrix} 0.62246 \\ 0.64566 \\ 0.66819 \end{bmatrix} \cdot \begin{pmatrix} 0.37754 \\ 0.38934 \\ 0.33181 \end{pmatrix} =$$

$$= \begin{bmatrix} -0.021794 \\ -0.05653 \\ 0.119323 \end{bmatrix}$$

$$\frac{\Delta E}{\Delta W^{[1]}} = \begin{bmatrix} -0.021794 \\ -0.05653 \\ 0.119323 \end{bmatrix} [2 \ 2] = \begin{bmatrix} -0.04359 & -0.04359 \\ -0.11306 & -0.11306 \\ -0.23865 & -0.23865 \end{bmatrix}$$

$$W^{[1]}_{\text{New}} = \begin{bmatrix} 0.1 & 0.1 \\ 0.1 & 0.2 \\ 0.2 & 0.1 \end{bmatrix} - 0.5 \begin{bmatrix} -0.04359 & -0.04359 \\ -0.11306 & -0.11306 \\ -0.23865 & -0.23865 \end{bmatrix}$$

$$= \begin{bmatrix} 0.1218 & 0.1218 \\ 0.18653 & 0.25653 \\ 0.31933 & 0.21933 \end{bmatrix}$$

To use sigmoid activation, allows the MLP to learn non-linear patterns and relationships, allowing for better prediction and representation capabilities.