

1.

D	y_1	y_2	y_{num}	y_{class}
x_1	1	1	1.25	B
x_2	1	3	7.0	A
x_3	3	2	2.7	C
x_4	3	3	3.2	A
x_5	2	4	5.5	B

$$w = (x^T \cdot x)^{-1} \cdot x^T \cdot z$$

Moore-Penrose
solution

$$X = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 3 \\ 1 & 3 & 2 \\ 1 & 3 & 3 \\ 1 & 2 & 4 \end{bmatrix} \quad z = \begin{bmatrix} 1.25 \\ 7.0 \\ 2.7 \\ 3.2 \\ 5.5 \end{bmatrix}$$

$$W = \left(\begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 3 & 3 & 2 \\ 1 & 3 & 2 & 3 & 4 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 3 \\ 1 & 3 & 2 \\ 1 & 3 & 3 \\ 1 & 2 & 4 \end{bmatrix} \right)^{-1} \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 3 & 3 & 2 \\ 1 & 3 & 2 & 3 & 4 \end{bmatrix} \begin{bmatrix} 1.25 \\ 7.0 \\ 2.7 \\ 3.2 \\ 5.5 \end{bmatrix} =$$

$$= \left(\begin{bmatrix} 5 & 10 & 13 \\ 10 & 24 & 27 \\ 13 & 27 & 39 \end{bmatrix} \right)^{-1} \cdot \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 3 & 3 & 2 \\ 1 & 3 & 2 & 3 & 4 \end{bmatrix} \begin{bmatrix} 1.25 \\ 7.0 \\ 2.7 \\ 3.2 \\ 5.5 \end{bmatrix} =$$

$$\begin{bmatrix} \frac{23}{11} & \frac{-13}{33} & \frac{-14}{33} \\ \frac{-13}{33} & \frac{26}{99} & \frac{-5}{99} \\ \frac{-14}{33} & \frac{-5}{99} & \frac{20}{99} \end{bmatrix} \cdot \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 3 & 3 & 2 \\ 1 & 3 & 2 & 3 & 4 \end{bmatrix} \begin{bmatrix} 1.25 \\ 7.0 \\ 2.7 \\ 3.2 \\ 5.5 \end{bmatrix} =$$

$$= \begin{bmatrix} \frac{14}{11} & \frac{14}{33} & \frac{2}{33} & \frac{-4}{11} & \frac{-13}{33} \\ \frac{-2}{11} & \frac{-28}{99} & \frac{29}{99} & \frac{8}{33} & \frac{-7}{99} \\ \frac{-3}{11} & \frac{13}{99} & \frac{-17}{99} & \frac{1}{33} & \frac{28}{99} \end{bmatrix} \begin{bmatrix} 1.25 \\ 7.0 \\ 2.7 \\ 3.2 \\ 5.5 \end{bmatrix} = \begin{bmatrix} \frac{46}{33} \\ \frac{-1019}{990} \\ \frac{3499}{1980} \end{bmatrix}$$

$$\text{Output}(y_1, y_2) = \frac{46}{33} - \frac{1019}{990} \cdot y_1 + \frac{3499}{1980} \cdot y_2$$

2. Ridge regression loss difference is in the loss function

$$E(w) = \frac{1}{2} \sum_{i=1}^N (z_i - \hat{z}_i)^2 + \frac{\lambda}{2} \|w\|^2$$

↓ Moore Penrose solution

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

$$X = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 3 \\ 1 & 3 & 2 \\ 1 & 3 & 3 \\ 1 & 2 & 4 \end{bmatrix} \quad z = \begin{bmatrix} 1.25 \\ 7.0 \\ 2.7 \\ 3.2 \\ 5.5 \end{bmatrix} \quad \lambda = 1$$

$$w = \left(\begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 3 & 3 & 2 \\ 1 & 3 & 2 & 3 & 4 \end{bmatrix} + \lambda \cdot I \right)^{-1} \cdot X^T \cdot z$$

$$= \left(\begin{bmatrix} 5 & 10 & 13 \\ 10 & 24 & 27 \\ 13 & 27 & 39 \end{bmatrix} + \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \right)^{-1} \cdot X^T \cdot z$$

$$= \left(\begin{bmatrix} 6 & 10 & 13 \\ 10 & 25 & 27 \\ 13 & 27 & 40 \end{bmatrix} \right)^{-1} \cdot X^T \cdot z =$$

$$= \begin{bmatrix} \frac{271}{421} & \frac{-49}{421} & \frac{55}{421} \\ \frac{-49}{421} & \frac{91}{421} & \frac{-32}{421} \\ \frac{55}{421} & \frac{-32}{421} & \frac{59}{421} \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 3 & 3 & 2 \\ 1 & 3 & 2 & 3 & 4 \end{bmatrix} \begin{bmatrix} 1.25 \\ 7.0 \\ 2.7 \\ 3.2 \\ 5.5 \end{bmatrix} =$$

$$= \begin{bmatrix} \frac{167}{421} & \frac{57}{421} & \frac{14}{421} & \frac{-41}{421} & \frac{-47}{421} \\ \frac{-10}{421} & \frac{-74}{421} & \frac{100}{421} & \frac{68}{421} & \frac{-35}{421} \\ \frac{-34}{421} & \frac{63}{421} & \frac{-51}{421} & \frac{-1}{421} & \frac{81}{421} \end{bmatrix} \begin{bmatrix} 1.25 \\ 7.0 \\ 2.7 \\ 3.2 \\ 5.5 \end{bmatrix} =$$

$$= \begin{bmatrix} \frac{5117}{8420} \\ \frac{-1177}{2105} \\ \frac{13987}{8420} \end{bmatrix} \quad w_1 = \begin{bmatrix} \frac{46}{33} \\ \frac{-1019}{990} \\ \frac{3499}{1980} \end{bmatrix} \quad w_2 = \begin{bmatrix} \frac{5117}{8420} \\ \frac{-1177}{2105} \\ \frac{13987}{8420} \end{bmatrix}$$

Now we compare the norm of weight vectors.
Ridge has a lower norm.

$$\|w_1\| = \frac{6136619}{71280} = 6.12541$$

$$\|w_2\| = \frac{121992561}{35448200} = 3.4416$$

Even if we remove the bias Ridge also has higher weight values.

$$\|w_1\| = \frac{3279284}{784080} = 4.18$$

$$\|w_2\| = \frac{217801433}{7086400} = 3.07$$

$$b = 1.3934$$

simple linear regression

$$b = 0.6077$$

Ridge regression

Both weights and bias are significantly penalized in Ridge regression.

$$3. \text{ RMSE} = \sqrt{\sum_{i=1}^n \frac{(z_i - \hat{z}_i)^2}{n}}$$

$$x_6 = \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix} \quad x_7 = \begin{bmatrix} 1 \\ 1 \\ 2 \end{bmatrix} \quad x_8 = \begin{bmatrix} 1 \\ 5 \\ 1 \end{bmatrix}$$

Simple linear regression:

$$w = \begin{bmatrix} \frac{46}{33} \\ \frac{-1019}{990} \\ \frac{3499}{1980} \end{bmatrix} \quad \hat{z}_6 = \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix} \cdot \begin{bmatrix} \frac{46}{33} \\ \frac{-1019}{990} \\ \frac{3499}{1980} \end{bmatrix} = \frac{947}{330} = 2.8697$$

$$\hat{z}_7 = \begin{bmatrix} 1 \\ 1 \\ 2 \end{bmatrix} \cdot \begin{bmatrix} \frac{46}{33} \\ \frac{-1019}{990} \\ \frac{3499}{1980} \end{bmatrix} = \frac{386}{99} = 3.8990$$

$$\hat{z}_8 = \begin{bmatrix} 1 \\ 5 \\ 1 \end{bmatrix} \cdot \begin{bmatrix} \frac{46}{33} \\ \frac{-1019}{990} \\ \frac{3499}{1980} \end{bmatrix} = \frac{-3931}{1980} = -1.9854$$

$$\text{RMSE} = \sqrt{\frac{(0.7 - 2.8697)^2 + (1.1 - 3.8990)^2 + (2.2 - (-1.9854))^2}{3}}$$

$$= \sqrt{\frac{4.7076 + 7.8344 + 4.1854}{3}} = 2.361$$

Ridge Regression

$$w = \begin{bmatrix} \frac{5117}{8420} \\ \frac{-1177}{2105} \\ \frac{13987}{8420} \end{bmatrix} \quad \hat{z}_6 = \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix} \cdot w = \frac{4735}{1684} = 2.8118$$

$$\hat{z}_7 = \begin{bmatrix} 1 \\ 1 \\ 2 \end{bmatrix} \cdot w = \frac{28383}{8420} = 3.3709$$

$$\hat{z}_8 = \begin{bmatrix} 1 \\ 5 \\ 1 \end{bmatrix} \cdot w = \frac{-1109}{2105} = -0.5268$$

$$\sqrt{\frac{(0.7 - 2.8118)^2 + (1.1 - 3.3709)^2 + (2.2 + 0.5268)^2}{3}} =$$

$$= \sqrt{\frac{4.61597}{3} + \frac{5.1570}{3} + \frac{7.43543824}{3}} =$$

$$= \sqrt{5.68410} = 2.384$$

Simple linear regression $\longrightarrow 2.361$

Ridge linear regression $\longrightarrow 2.384$

Ridge regression usually prevents overfitting. In this case it did not improve the model's accuracy so it is likely that linear regression was underfitting, that is why ridge worsened the results.

$$4. \quad W^1 = \begin{bmatrix} 0.1 & 0.1 \\ 0.1 & 0.2 \\ 0.2 & 0.1 \end{bmatrix} \quad b^1 = \begin{bmatrix} 0.1 \\ 0 \\ 0.1 \end{bmatrix}$$

$$W^2 = \begin{bmatrix} 1 & 2 & 2 \\ 1 & 2 & 1 \\ 1 & 1 & 1 \end{bmatrix} \quad b^2 = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

$$x_1 = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \quad W^1 x_1 + b^1 = \begin{bmatrix} 0.1 & 0.1 \\ 0.1 & 0.2 \\ 0.2 & 0.1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 0.1 \\ 0 \\ 0.1 \end{bmatrix} =$$

$$= \begin{bmatrix} 0.2 \\ 0.3 \\ 0.3 \end{bmatrix} + \begin{bmatrix} 0.1 \\ 0 \\ 0.1 \end{bmatrix} = \begin{bmatrix} 0.3 \\ 0.3 \\ 0.4 \end{bmatrix} \rightarrow z^1$$

$$W^2 z^1 + b^2 = \begin{bmatrix} 1 & 2 & 2 \\ 1 & 2 & 1 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 0.3 \\ 0.3 \\ 0.4 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} =$$

$$= \begin{bmatrix} 14 \\ 13 \\ 10 \\ 1 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 20 \\ 13 \\ 10 \\ 2 \end{bmatrix} \rightarrow z^2$$

$$\sum_{i=1}^N e^{z_i^2} = e^{\frac{20}{3}} + e^{\frac{23}{10}} + e^2 = 803.1352$$

$$\frac{e^{\frac{20}{3}}}{803.1352} = 0.97838 \quad \hat{z}_1 = 4$$

$$\frac{e^{\frac{23}{10}}}{803.1352} = 0.01241 \quad \begin{bmatrix} 0.97838 \\ 0.01241 \\ 0.00920 \end{bmatrix}$$

$$\frac{e^2}{803.1352} = 0.00920$$

$$\text{Our targets for } x_1 \text{ are: } z_1 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

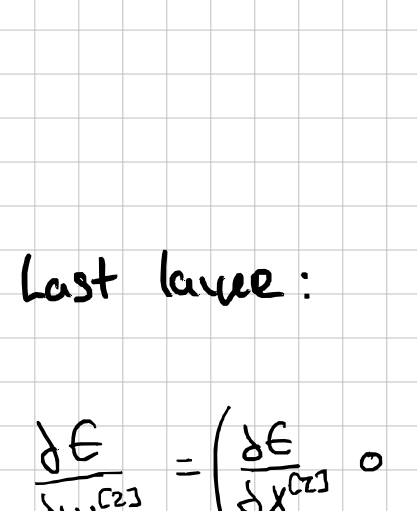
To perform a stochastic gradient descent, we have to optimize the following function:

$$- \sum_{i=1}^N \sum_{l=1}^{L-1} t_l^i \log(z_l^{\text{out}, i}) =$$

$$= - (t_2 \log z_2^{\text{out}}) = - \log z_2^{\text{out}}$$

$$t_2 = 1$$

$$- \log(\text{softmax}(x_2^{\text{out}}))$$



Stochastic Gradient descent for x_1

$$W_{\text{new}}^{[L]} = W_{\text{old}}^{[L]} - \mu \frac{\partial E}{\partial W^{[L]}}$$

$$\delta^{[L]} = x^{[L]} - t \quad \text{(last layer)}$$

$$\frac{\partial E}{\partial W^{[L]}} = \underbrace{\left(\frac{\partial E}{\partial x^{[L]}} \circ \frac{\partial x^{[L]}}{\partial z^{[L]}} \right)}_{\delta^{[L]}} \cdot \left(\frac{\partial z^{[L]}}{\partial W^{[L]}} \right)^T$$

Last layer:

$$\frac{\partial E}{\partial W^{[L]}} = \underbrace{\left(\frac{\partial E}{\partial x^{[L]}} \circ \frac{\partial x^{[L]}}{\partial z^{[L]}} \right)}_{\delta^{[L]}} \cdot \underbrace{\left(\frac{\partial z^{[L]}}{\partial W^{[L]}} \right)^T}_{x^{[L]}} =$$

$$= \delta^{[L]} (x^{[L]})^T$$

$$\delta^{[L]} = x^{[L]} - t = \begin{bmatrix} 0.97838 \\ 0.01241 \\ 0.00920 \end{bmatrix} - \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0.97838 \\ -0.98759 \\ 0.00920 \end{bmatrix}$$

$$\frac{\partial E}{\partial W^{[L]}} = \begin{bmatrix} 0.97838 \\ -0.98759 \\ 0.00920 \end{bmatrix} \begin{bmatrix} 0.3 & 0.3 & 0.4 \end{bmatrix} =$$

$$= \begin{bmatrix} 0.2935 & 0.2935 & 0.3914 \\ -0.2963 & -0.2963 & -0.3950 \\ 0.0028 & 0.0028 & 0.0037 \end{bmatrix}$$

$$W_{\text{new}}^{[L]} = W_{\text{old}}^{[L]} - \mu \frac{\partial E}{\partial W^{[L]}} =$$

$$= \begin{bmatrix} 1 & 2 & 2 \\ 1 & 2 & 1 \\ 1 & 1 & 1 \end{bmatrix} - 0.1 \begin{bmatrix} 0.2935 & 0.2935 & 0.3914 \\ -0.2963 & -0.2963 & -0.3950 \\ 0.0028 & 0.0028 & 0.0037 \end{bmatrix}$$

$$= \begin{bmatrix} 0.9706 & 1.9706 & 1.9609 \\ 1.0296 & 2.0296 & 1.0395 \\ 0.9997 & 0.9997 & 0.9996 \end{bmatrix}$$

Remaining layer

$$\delta^{[L]} = \left(\frac{\partial z^{[L+1]}}{\partial x^{[L]}} \right)^T \cdot \delta^{[L+1]} \circ \frac{\partial x^{[L]}}{\partial z^{[L]}}$$

$$= W^{[L+1]T} \cdot \delta^{[L+1]} \circ \frac{\partial x^{[L]}}{\partial z^{[L]}} = W^{[L+1]T} \cdot \delta^{[L+1]} \quad (\text{no activation function})$$

$$\delta^{[L]} = W^{[L+1]T} \cdot \delta^{[L+1]} =$$

$$= \begin{bmatrix} 0.9706 & 1.0296 & 0.9997 \\ 1.9706 & 2.0296 & 0.9997 \\ 1.9609 & 1.0395 & 0.9996 \end{bmatrix} \begin{bmatrix} 0.97838 \\ -0.98759 \\ 0.00920 \end{bmatrix}$$

$$= \begin{bmatrix} -0.0580 \\ -0.0672 \\ 0.9011 \end{bmatrix}$$

$$\frac{\partial E}{\partial W^{[L]}} = \delta^{[L]} \cdot (x^{[L-1]})^T$$

$$\frac{\partial E}{\partial W^{[L]}} = \delta^{[L]} \cdot (x^{[0]})^T =$$

$$= \begin{bmatrix} -0.0580 \\ -0.0672 \\ 0.9011 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1.25 \end{bmatrix} = \begin{bmatrix} -0.0580 & -0.0580 \\ -0.0672 & -0.0672 \\ 0.9011 & 0.9011 \end{bmatrix}$$

$$W_{\text{new}}^{[L]} = W_{\text{old}}^{[L]} - \mu \frac{\partial E}{\partial W^{[L]}} =$$

$$= \begin{bmatrix} 0.1 & 0.1 \\ 0.1 & 0.2 \\ 0.2 & 0.1 \end{bmatrix} - \begin{bmatrix} -0.0580 & -0.0580 \\ -0.0672 & -0.0672 \\ 0.9011 & 0.9011 \end{bmatrix} =$$

$$= \begin{bmatrix} 0.1580 & 0.1580 \\ 0.1672 & 0.2672 \\ -0.7011 & -0.8011 \end{bmatrix}$$

Bias update:

$$\frac{\partial E}{\partial b^{[L]}} = \delta^{[L]} \cdot \frac{\partial z^{[L]}}{\partial b^{[L]}} = \delta^{[L]}$$

$$b_{\text{new}}^{[L]} = b_{\text{old}}^{[L]} - \mu \frac{\partial E}{\partial b^{[L]}} = b_{\text{old}}^{[L]} - \mu \delta^{[L]}$$

$$= \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} - 0.1 \begin{bmatrix} 0.97838 \\ -0.98759 \\ 0.00920 \end{bmatrix} = \begin{bmatrix} 0.9022 \\ 1.9088 \\ 0.9991 \end{bmatrix}$$

$$b_{\text{new}}^{[1]} = b_{\text{old}}^{[1]} - \mu \frac{\partial E}{\partial b^{[1]}} = b_{\text{old}}^{[1]} - \mu \delta^{[1]}$$

$$= \begin{bmatrix} 0.1 \\ 0 \\ 0.1 \end{bmatrix} - 0.1 \begin{bmatrix} -0.0580 \\ -0.0672 \\ 0.9011 \end{bmatrix} = \begin{bmatrix} 0.1058 \\ 0.0067 \\ 0.0099 \end{bmatrix}$$

Summary of updates in ex. 4

$$W^{[1]} = \begin{bmatrix} 0.1580 & 0.1580 \\ 0.1672 & 0.2672 \\ -0.7011 & -0.8011 \end{bmatrix}$$

$$b^{[1]} = \begin{bmatrix} 0.1058 \\ 0.0067 \\ 0.0099 \end{bmatrix}$$

$$b^{[2]} = \begin{bmatrix} 0.9022 \\ 1.9088 \\ 0.9991 \end{bmatrix}$$

$$W^{[2]} = \begin{bmatrix} 0.9706 & 1.9706 & 1.9609 \\ 1.0296 & 2.0296 & 1.0395 \\ 0.9997 & 0.9997 & 0.9996 \end{bmatrix}$$