

Proof of gradient update rule  
 prove that  $\delta_{ij}$  is equivalent to  $\frac{\partial L}{\partial z_{ij}}$  for  $0 \leq i \leq L$

$$L = \frac{1}{2} (y - \hat{y})^2; \quad \frac{\partial L}{\partial z_{L0}} = \frac{\partial L}{\partial \hat{y}} = \hat{y} - y$$

$$z_{L0} = \sum_{p=0}^{h-1} w_{yp} \phi_{Lp} + b_{\hat{y}} = \hat{y}$$

$$\frac{\partial L}{\partial z_{(L-1)j}} = \frac{\partial L}{\partial z_{L0}} \cdot \frac{\partial z_{L0}}{\partial \phi_{Lj}} \cdot \frac{\partial \phi_{Lj}}{\partial z_{(L-1)j}} = \frac{\partial L}{\partial z_{L0}} \cdot w_{\hat{y}j} \cdot A'(z_{(L-1)j})$$

$$\frac{\partial L}{\partial z_{(L-2)j}} = \sum_{p=0}^{h-1} \frac{\partial L}{\partial z_{(L-1)p}} \cdot \frac{\partial z_{(L-1)p}}{\partial \phi_{(L-1)j}} \cdot \frac{\partial \phi_{(L-1)j}}{\partial z_{(L-2)j}}$$

$$= \sum_{p=0}^{h-1} \frac{\partial L}{\partial z_{(L-1)p}} \cdot w_{\phi_{(L-1)p}j} \cdot A'(z_{(L-2)j})$$

Statement  $\delta_{ij}^c \frac{\partial L}{\partial z_{ij}} = \sum_{p=0}^{h-1} \frac{\partial L}{\partial z_{(c+1)p}} \cdot \frac{\partial z_{(c+1)p}}{\partial \phi_{(c+1)j}} \cdot \frac{\partial \phi_{(c+1)j}}{\partial z_{ij}}$

$\delta_{L-2}$  is true, show  $\delta_{i+1} \rightarrow \delta_i$

$$\begin{aligned} \frac{\partial L}{\partial z_{ij}} &= \sum_{p=0}^{h-1} \left( \sum_{k=0}^{h-1} \frac{\partial L}{\partial z_{(i+2)k}} \cdot \frac{\partial z_{(i+2)k}}{\partial \phi_{(i+2)p}} \cdot \frac{\partial \phi_{(i+2)p}}{\partial z_{(i+1)p}} \right) \cdot \frac{\partial z_{(i+1)p}}{\partial \phi_{(i+1)j}} \cdot \frac{\partial \phi_{(i+1)j}}{\partial z_{ij}} \\ &= \sum_{p=0}^{h-1} \frac{\partial L}{\partial z_{(i+1)p}} \cdot \frac{\partial z_{(i+1)p}}{\partial \phi_{(i+1)j}} \cdot \frac{\partial \phi_{(i+1)j}}{\partial z_{ij}} \end{aligned}$$

$$\frac{\partial L}{\partial z_{ij}} = \sum_{p=0}^{h-1} \frac{\partial L}{\partial z_{(i+1)p}} \cdot w_{\phi_{(i+1)p}j} \cdot A'(z_{ij}) \quad \text{for } 0 \leq i \leq L-2$$

So we can represent  $\frac{\partial L}{\partial z_{ij}}$  as  $\delta_{ij}$

~~$$\delta_{L0} = \hat{y} - y$$~~

$$\delta_{(L-1)j} = (\hat{y} - y) \cdot w_{\hat{y}j} \cdot A'(z_{(L-1)j})$$

$$\delta_{L0} = \hat{y} - y \quad \delta_{(L-2)j} = \delta_{ij} = \sum_{p=0}^{h-1} \delta_{(i+1)p} \cdot w_{\phi_{(i+1)p}j} \cdot A'(z_{ij}) \quad \forall 0 \leq i \leq L-2$$

backward pass for some  $x_c, y_c$   
error backpropagation

$$\delta = \hat{y} - y; \quad T(\delta_L) = \theta(1)$$

error term for neuron  $i_j$  ( $L_0$  is output neuron)

$$\delta_{(L-1)j} = \delta_{L_0} \cdot w_{g_j} \cdot A'(z_{(L-1)j}); \quad T(\delta_{(L-1)j}) = \theta(h)$$

$$\delta_{(L-2)i} = \sum_{p=0}^{h-1} \delta_{(L-1)p} \cdot w_{\phi_{(L-1)pj}} \cdot A'(z_{(L-2)i}); \quad T(\delta_{(L-2)i}) = \theta(h^2)$$

$$\delta_{ij} = \sum_{p=0}^{h-1} \delta_{(L-1)p} \cdot w_{\phi_{(L-1)pj}} \cdot A'(z_{ij}); \quad T(\delta_{ij}) = \theta(h^2)$$

$$T(\delta) = \theta(1) + \theta(h) + \theta((L-2)h^2) \rightarrow T(\delta) = \theta(Lh^2)$$

gradient update

$$w_{\phi_{ijp}} = w_{\phi_{ijp}} - \eta \delta_{ij} \phi_{ip}$$

$$w_{g_p} = w_{g_p} - \eta \delta_{L_0} \phi_{Lp}$$

$$b_{\phi_{ij}} = b_{\phi_{ij}} - \eta \delta_{ij}$$

$$b_g = b_g - \eta \delta_{L_0}$$

$$T(\text{update}) = \theta(h) + \theta((L-2)h^2) + \theta(nh)$$

$$\rightarrow T(\text{update}) = \begin{cases} \theta(Lh^2) & \text{if } Lh > n \\ \theta(nh) & \text{else} \end{cases}$$

gradient clipping (done before update)

$$\delta_{ij} = \begin{cases} \psi \cdot \text{sign}(\delta_{ij}) & \text{if } |\delta_{ij}| > \psi \\ \delta_{ij} & \text{else} \end{cases} \quad \text{for some clipping threshold } \psi$$

$$T(\text{clipping}) = \theta(1) + \theta((L-1) \cdot h) = \theta(Lh)$$

$$T(\text{backward}) = T(\delta) + T(\text{clipping}) + T(\text{update})$$

$$\rightarrow T(\text{backward}) = \begin{cases} \theta(Lh^2) & \text{if } Lh > n \\ \theta(nh) & \text{else} \end{cases}$$

$$T(\text{pass}) = T(\text{forward}) + T(\text{backward})$$

$$\rightarrow T(\text{pass}) = \begin{cases} \theta(Lh^2) & \text{if } Lh > n \\ \theta(nh) & \text{else} \end{cases}$$

$$T(\text{pass}_{\text{total}}) = \begin{cases} \theta(MLh^2) & \text{if } Lh > n \\ \theta(Mnh) & \text{else} \end{cases}$$

for  $M$  samples

$$T(\text{train}) = \begin{cases} \theta(KMLh^2) & \text{if } Lh > n \\ \theta(KMnh) & \text{else} \end{cases}$$

for  $K$  training iterations

total time complexity

# Multilayer perceptron

Input/target

$$X = \begin{bmatrix} \begin{bmatrix} x_{c1} \\ \vdots \\ x_{cn} \end{bmatrix} \end{bmatrix}_{c=0}^{M-1}; Y = \begin{bmatrix} y_c \end{bmatrix}_{c=0}^{M-1}; M = \# \text{ samples in dataset}, n = \# \text{ input features}$$

$$T(\text{read}) = \theta(Mn) + \theta(M) = \theta(Mn)$$

other parameters:  $h$  = hidden layer size,  $L$  = hidden layer dim (# HLS)

$\eta$  = learning rate

Weight/bias initialization

$$SD_i = \begin{cases} \sqrt{\frac{2}{n}} & i=0 \\ \sqrt{\frac{2}{h}} & i>0 \end{cases} \quad ND_i = N(0, SD_i)$$

$zer = \text{random \# generator}$

$$W_\phi = \begin{bmatrix} \begin{bmatrix} ND_i(zer) \end{bmatrix}_{i=0}^{L-1} \\ \begin{bmatrix} ND_i(zer) \end{bmatrix}_{i=0}^{L-1} \end{bmatrix}_{j=0}^{n-1}$$

for  $0 \leq i \leq L$

$$W_\psi = \begin{bmatrix} ND_i(zer) \end{bmatrix}_{j=0}^{h-1} \quad b_\psi = ND_i(zer) \quad b_\phi = \begin{bmatrix} ND_i(zer) \end{bmatrix}_{i=0}^{L-1}$$

$$T(\text{init}) = \theta(1) + \theta(h) + \theta(Lh) + \theta((L-1)h^2) + \theta(nh)$$

$$\rightarrow T(\text{init}) = \begin{cases} \theta(Lh^2) & \text{if } Lh > n \\ \theta(nh) & \text{else} \end{cases}$$

Forward pass

for some  $x_c$

$$\phi = \begin{bmatrix} \begin{bmatrix} \phi_{ij} \end{bmatrix}_{i=0}^{L-1} \end{bmatrix}_{j=0}^{h-1}; \phi_{ij} = A(z_{ij})$$

for some activation function  $A$

$$T(A) = \theta(1)$$

$$z_{ij} = \begin{cases} \sum_{p=0}^{n-1} W_{\phi,ij} x_p + b_{\phi,ij} & i=0 \\ \sum_{p=0}^{h-1} W_{\phi,ij} \phi_{p,i-1} + b_{\phi,ij} & i>0 \end{cases}$$

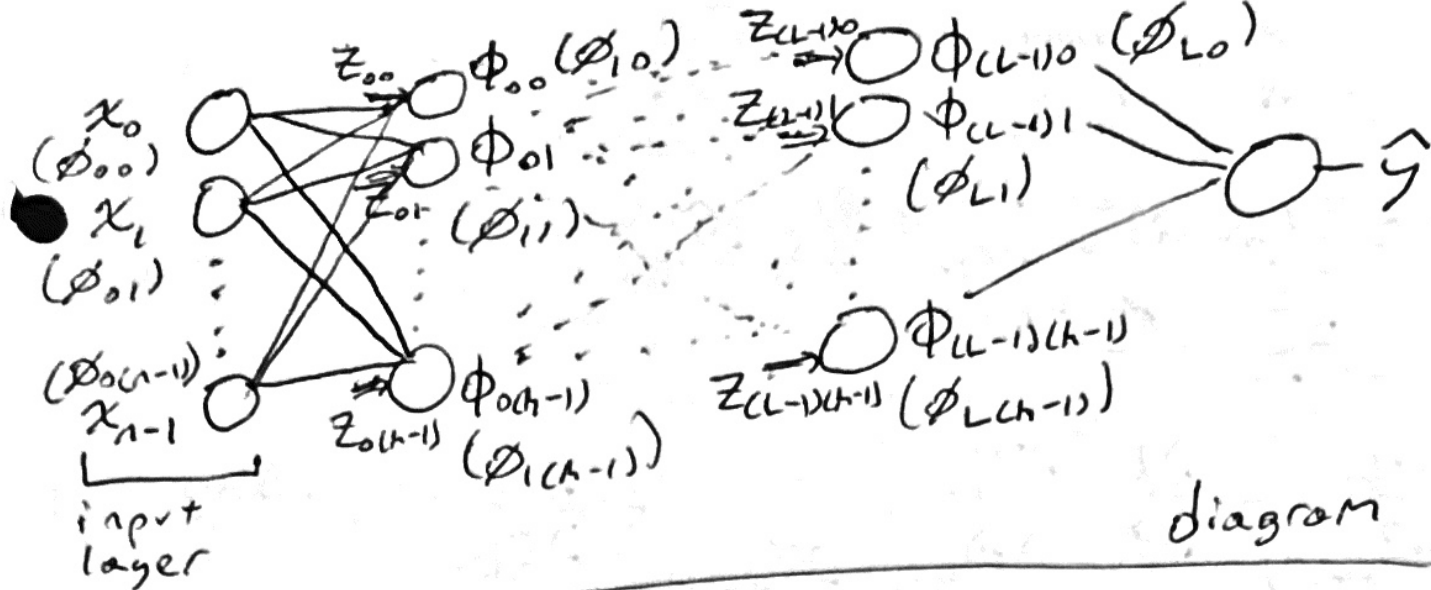
$$T(z_{ij}) = \begin{cases} \theta(1) & \text{if } i=0 \\ \theta(h) & \text{if } i>0 \end{cases}$$

$$T(\phi) = \theta((L-1)h^2) + \theta(nh) \quad T(\phi_{ij}) = T(z_{ij})$$

$$\rightarrow T(\phi) = \begin{cases} \theta(Lh^2) & \text{if } Lh > n \\ \theta(nh) & \text{else} \end{cases} \quad \hat{y} = \sum_{p=0}^{h-1} W_{\psi,p} \phi_{p,L-1} + b_{\hat{y}}, \quad T(\hat{y}) = \theta(h)$$

$$\phi = \begin{bmatrix} \begin{bmatrix} x_c \end{bmatrix}_{i=0}^{n-1} \\ \begin{bmatrix} \phi_{c(i-1),j} \end{bmatrix}_{j=0}^{h-1} \end{bmatrix}_{i=0}^{L-1}; \phi = \begin{bmatrix} \phi_c \end{bmatrix}_{c=0}^{M-1}$$

$$T(\text{forward}) = T(\phi) + T(\hat{y}) = \begin{cases} \theta(Lh^2) & \text{if } Lh > n \\ \theta(nh) & \text{else} \end{cases}$$



activation functions  
 $\text{relu}(z) = \begin{cases} z & \text{if } z > 0 \\ 0 & \text{else} \end{cases}; \text{relu}'(z) = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{else} \end{cases}$

$\text{leaky relu}(z) = \begin{cases} z & \text{if } z > 0 \\ \alpha z & \text{else} \end{cases}; \text{leaky relu}'(z) = \begin{cases} 1 & \text{if } z > 0 \\ \alpha & \text{else} \end{cases}$   
 $\alpha$  is a small constant

$\sigma(z) = \frac{1}{1+e^{-z}}; \sigma'(z) = \sigma(z)(1-\sigma(z))$

$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}; \tanh'(z) = 1 - \tanh^2(z)$

$\forall A \text{ above, } T(A) = \theta(1), T(A') = \theta(1)$

time recordings  
 $M = 15000 = 15 \cdot 10^3$   
 $n = 6528 \quad n > Lh$   
 $h = 5 \rightarrow T(\text{train}) = \theta(kmh)$   
 $L = 2 \quad = C \cdot 3 \cdot 6528 \cdot 5 = C \cdot 97920 \cdot 15 \cdot 10^3$   
 $k = 3 \quad = C \cdot 1,468,800 \cdot 10^3$   
 running time!  
 199 seconds  
 $= 199 \cdot 10^9 \text{ ns}$   
 running time!  
 189 seconds  
 $C = \frac{189}{1,468,800} \approx 129$   
 $C = \frac{199 \cdot 10^9}{1,468,800 \cdot 10^9} \approx 132.08$