

$$E(X) = \frac{1}{2} \sum_t \sum_m (r_m^t - O_m^{(3)t})^2$$

$$w_{t+1} = w_t + \eta \cdot \frac{\partial E}{\partial w}$$

$$\Delta w = -\eta \frac{\partial E}{\partial w}$$

$$\frac{\partial E}{\partial a_{k,m}^{(2)}} = \frac{\partial}{\partial a_{k,m}^{(2)}} \left( \frac{1}{2} \sum_t \sum_m (r_m^t - O_m^{(3)t})^2 \right)$$

$$\Delta a_{k,m}^{(2)} = -\eta \cdot \frac{\partial E}{\partial a_{k,m}^{(2)}}$$

$$\begin{aligned} &= \sum_t \sum_m (r_m^t - O_m^{(3)t}) \cdot \frac{\partial O_m^{(3)t}}{\partial a_{k,m}^{(2)}} \\ &= X_k^{(2)t} \end{aligned}$$

$$\Delta a_{k,m}^{(2)} = \sum_t (r_m^t - O_m^{(3)t}) \cdot X_k^{(2)t}$$

$$a^{(0)} = \begin{bmatrix} a_{0,1}^{(0)} & a_{0,2}^{(0)} & a_{0,3}^{(0)} & a_{0,4}^{(0)} \\ a_{1,1}^{(0)} & a_{1,2}^{(0)} & a_{1,3}^{(0)} & a_{1,4}^{(0)} \\ a_{2,1}^{(0)} & a_{2,2}^{(0)} & a_{2,3}^{(0)} & a_{2,4}^{(0)} \\ a_{3,1}^{(0)} & a_{3,2}^{(0)} & a_{3,3}^{(0)} & a_{3,4}^{(0)} \end{bmatrix}$$

$$\begin{aligned} X_0^{(l)} &= X_0^{(1)} = X_0^{(2)} = +1 \\ X_1^{(l)} &= [X_0^{(l)} \ X_1^{(l)} \ X_2^{(l)} \ X_3^{(l)} \dots] \\ O^{(l)} &= [+1 \ O_1^{(l)} \ O_2^{(l)} \ O_3^{(l)} \dots] \end{aligned}$$

added after

$$O_s^{(l)t} = \sum_a a_{a,s}^{(l-1)t} \cdot X_a^{(l-1)t} \rightarrow \text{expand } a \rightarrow =$$

$$= \begin{bmatrix} a_{0,0}^{(l-1)} & a_{1,0}^{(l-1)} & a_{2,0}^{(l-1)} & a_{3,0}^{(l-1)} \\ a_{0,1}^{(l-1)} & a_{1,1}^{(l-1)} & a_{2,1}^{(l-1)} & a_{3,1}^{(l-1)} \\ a_{0,2}^{(l-1)} & a_{1,2}^{(l-1)} & a_{2,2}^{(l-1)} & a_{3,2}^{(l-1)} \\ a_{0,3}^{(l-1)} & a_{1,3}^{(l-1)} & a_{2,3}^{(l-1)} & a_{3,3}^{(l-1)} \end{bmatrix} \begin{bmatrix} X_0 \ X_1 \ X_2 \ X_3 \end{bmatrix}$$

scalar

$$\begin{bmatrix} a_{0,5}^{(l-1)} & a_{1,5}^{(l-1)} & a_{2,5}^{(l-1)} & a_{3,5}^{(l-1)} \end{bmatrix} \cdot [X_0 \ X_1 \ X_2 \ X_3]$$

matrix multiplication

$$O^{(l)t} = X^{(l-1)t} \cdot a^{(l-1)}$$

matrix form

$$\begin{aligned} &[X_0 \ X_1 \ X_2 \ X_3] \cdot \begin{bmatrix} a_{0,0}^{(l-1)} & a_{0,1}^{(l-1)} & a_{0,2}^{(l-1)} \\ a_{1,0}^{(l-1)} & a_{1,1}^{(l-1)} & a_{1,2}^{(l-1)} \\ a_{2,0}^{(l-1)} & a_{2,1}^{(l-1)} & a_{2,2}^{(l-1)} \\ a_{3,0}^{(l-1)} & a_{3,1}^{(l-1)} & a_{3,2}^{(l-1)} \end{bmatrix} \\ &= [+1 \ O_1^{(l)} \ O_2^{(l)} \ O_3^{(l)} \dots] \end{aligned}$$

$$2.) \Delta a_{k,m}^{(1)} = \sum_t \sum_m (r_m^t - O_m^{(1)t}) \cdot X_k^{(1)t}$$

$$= \sum_t \sum_m e_m^t \cdot X_k^{(2)t} \rightarrow \text{open m index} \rightarrow \sum_t [e_0^t e_1^t] \cdot X_k^{(2)t}$$

$$\Delta a_k^{(2)} = \sum_t [e_0^t e_1^t] \cdot X_k^{(2)t}$$

$$\Delta a^{(2)} = \left[ \begin{array}{l} \sum_t [e_0^t e_1^t] \cdot X_0^{(2)t} \\ \sum_t [e_0^t e_1^t] \cdot X_1^{(2)t} \\ \sum_t [e_0^t e_1^t] \cdot X_2^{(2)t} \\ \vdots \end{array} \right] \xrightarrow{\text{matrix multiplication}} \left[ \begin{array}{l} X_0^{(2)t} \\ X_1^{(2)t} \\ X_2^{(2)t} \\ \vdots \end{array} \right] \cdot [e_0^t e_1^t] \rightarrow X^{(2)t} \cdot err$$

$$\Delta a_{j,k}^{(1)} = -\eta \frac{\partial E}{\partial a_{j,k}^{(1)}} \quad \frac{\partial E}{\partial a_{j,k}^{(1)}} = \partial \left( \frac{1}{2} \sum_t \sum_m (r_m^t - O_m^{(1)t})^2 \right) \quad err^t = [e_0^t e_1^t]$$

$$= \sum_t \sum_m (r_m^t - O_m^{(1)t})^2 - 1 \cdot \frac{\partial O_m^{(1)t}}{\partial a_{j,k}^{(1)}}$$

$$= a_{k,m}^{(2)} \cdot \frac{\partial X_k^{(2)t}}{\partial a_{j,k}^{(1)}} \quad \frac{\partial X_k^{(2)t}}{\partial a_{j,k}^{(1)}} = \frac{\partial (\sigma(O_k^{(2)t}))}{\partial a_{j,k}^{(1)}} \quad \frac{\partial O_m^{(1)t}}{\partial a_{j,k}^{(1)}} = \frac{\partial \left( \sum_{k=1}^n a_{k,m}^{(2)} \cdot X_k^{(2)t} \right)}{\partial a_{j,k}^{(1)}}$$

$$\frac{\partial O_k^{(2)t}}{\partial a_{j,k}^{(1)}} = \frac{\partial \left( \sum_{j=1}^m a_{j,k}^{(1)} \cdot X_j^{(1)t} \right)}{\partial a_{j,k}^{(1)}} = X_j^{(1)t}$$

$$\Delta a_{j,k}^{(1)} = +\eta \sum_t \sum_m e_m^t + 1 \cdot a_{k,m}^{(2)} \cdot \sigma(O_k^{(2)t}) \cdot (1 - \sigma(O_k^{(2)t})) \cdot X_j^{(1)t}$$

$$\sigma(O_k^{(2)t}) = \sigma(O_k^{(2)t})(1 - \sigma(O_k^{(2)t}))$$

$$\Delta a_{j,k}^{(1)} = \eta \sum_t \sum_m e_m^t \cdot a_{k,m}^{(2)} \cdot \sigma(O_k^{(2)t}) \cdot X_j^{(1)t}$$

→ open m in vector form

$$\Delta a_{j,k}^{(1)} = \eta \cdot \sum_t [e_0^t e_1^t] [a_{k,1}^{(2)}, a_{k,2}^{(2)}, \dots, a_{k,5}^{(2)}] \cdot \sigma(O_k^{(2)t}) \cdot X_j^{(1)t}$$

→  $[a_{0,1}^{(2)}, a_{0,2}^{(2)}, \dots, a_{0,5}^{(2)}] \rightarrow \text{excluded!}$

$$\begin{bmatrix} a_{1,1}^{(2)} & a_{1,2}^{(2)} \\ a_{2,1}^{(2)} & a_{2,2}^{(2)} \\ a_{3,1}^{(2)} & a_{3,2}^{(2)} \\ a_{4,1}^{(2)} & a_{4,2}^{(2)} \\ a_{5,1}^{(2)} & a_{5,2}^{(2)} \end{bmatrix} \cdot \begin{bmatrix} \sigma(O_1^{(2)t}) & \sigma(O_2^{(2)t}) & \sigma(O_3^{(2)t}) & \sigma(O_4^{(2)t}) & \sigma(O_5^{(2)t}) \end{bmatrix} \cdot X_j^{(1)t}$$

→  $a_{j,k}^{(1)} \rightarrow 5 \times 5 \text{ matrix}$   
 j is the row number  
 k is the column number  
 should be 5x5

→ This part should be smaller  
 dimensions 5x5

$$\eta \cdot \sum_t (err^t \cdot a^{(2)t}) \odot [\sigma(O_1^{(2)t}) \sigma(O_2^{(2)t}) \dots \sigma(O_5^{(2)t})]^T$$

→ pairwise  
 matrix multiplication

$$\sum_t err^t \cdot a^{(2)t} \odot \sigma(O^{(2)t}) \cdot X_j^{(1)t}$$

$$\Delta a_{j,k}^{(1)} = \sum_t err^t \cdot a^{(2)t} \odot \sigma(O^{(2)t}) \cdot X_j^{(1)t}$$

$$\Delta a_{j,k}^{(1)} = \gamma \sum_m e_m^t \cdot a_{k,m}^{(2)} \alpha(O_k^{(2)t}) \cdot x_j^{(1)t}$$

↳ vectorize with m

$$\Delta a_{j,k}^{(1)} = \gamma \sum_t [e_0^t \ e_1^t] [a_{k,0}^{(2)} \ a_{k,1}^{(2)}] \alpha(O_k^{(2)t}) \cdot x_j^{(1)t}$$

$[a_{0,0}^{(2)} \ a_{0,1}^{(2)}] \rightarrow \text{excluded}$

$$\Delta a_j^{(1)} = \gamma \sum_t [e_0^t \ e_1^t] \cdot \begin{bmatrix} a_{1,0}^{(2)} & a_{1,1}^{(2)} \\ a_{2,0}^{(2)} & a_{2,1}^{(2)} \\ a_{3,0}^{(2)} & a_{3,1}^{(2)} \\ a_{4,0}^{(2)} & a_{4,1}^{(2)} \\ a_{5,0}^{(2)} & a_{5,1}^{(2)} \end{bmatrix} [\alpha(O_1^{(2)t}) \alpha(O_2^{(2)t}) \dots \alpha(O_5^{(2)t})] \cdot x_j^{(1)t}$$

↳ whole expression here must be a row column

this is a vector  
and row dimension is 5

summation could be obtained  
 $[e_0^t \ e_1^t]$  pairwise could be multiplied to  $a^{(2)}$   
but error could be an matrix like

$$\begin{bmatrix} e_0^0 & e_1^0 \\ e_0^1 & e_1^1 \\ e_0^2 & e_1^2 \end{bmatrix} \rightarrow \text{so pairwise multiplication could not be used}$$

So matrix multiplication should be used.

$a^{(2)} \cdot \text{err}^T$  could be done

$$\begin{bmatrix} a_{1,0}^{(2)} & a_{1,1}^{(2)} \\ a_{2,0}^{(2)} & a_{2,1}^{(2)} \\ a_{3,0}^{(2)} & a_{3,1}^{(2)} \end{bmatrix} \begin{bmatrix} e_0^0 & e_0^1 \\ e_1^0 & e_1^1 \end{bmatrix} \rightarrow \text{this product a matrix from example}$$

$$\begin{bmatrix} \quad \end{bmatrix}_{5 \times \text{sample number}}$$

but  $\alpha(O_k^{(2)})$   
is in dimens

$5 \times \text{sample count}$   
So to multiply b

extra transpose operations should be performed

Instead let's look err.  $a^{(2)T} [e_0^t \ e_1^t] \begin{bmatrix} a_{1,0} & a_{2,0} & a_{3,0} & a_{4,0} & a_{5,0} \\ a_{1,1} & a_{2,1} & a_{3,1} & a_{4,1} & a_{5,1} \end{bmatrix}$

↳ this produces  $5 \times \text{sample count}$   
 $\alpha(O^{(2)t})$  is also  $5 \times \text{sample count}$  therefore pairwise

$$\Delta a_j^{(1)} = \text{err. } a^{(2)T} \odot \alpha(O^{(2)t}) \cdot x_j^{(1)t} \text{ could be done.}$$

$$\Delta a_j^{(1)} = \left[ (e_0^t a_{1,0}^{(2)} + e_1^t a_{1,1}^{(2)}) \cdot \alpha(O_1^{(2)t}) \quad (e_0^t a_{2,0}^{(2)} + e_1^t a_{2,1}^{(2)}) \cdot \alpha(O_2^{(2)t}) \quad \dots \right]$$

$$\Delta a^{(1)} = \left[ \begin{array}{c} \text{err. } a^{(2)T} \odot \alpha(O^{(2)t}) \cdot x_0^{(1)t} \\ \text{err. } a^{(2)T} \odot \alpha(O^{(2)t}) \cdot x_1^{(1)t} \\ \text{err. } a^{(2)T} \odot \alpha(O^{(2)t}) \cdot x_2^{(1)t} \end{array} \right]$$

$$\left[ \begin{array}{c} x_0^{(1)t} \\ x_1^{(1)t} \\ x_2^{(1)t} \end{array} \right] \left[ (\text{err. } a^{(2)T}) \odot \alpha(O^{(2)t}) \right]$$

$$\Delta a^{(1)} = \sum_t x^{(1)t} \cdot ((\text{err. } a^{(2)T}) \odot \alpha(O^{(2)t}))$$

$$\Delta a_{j,k}^{(1)} = \eta \sum_t \sum_m e_m^t \cdot a_{k,m}^{(2)} \alpha(O_k^{(2)t}) \cdot X_j^{(1)t}$$

$$\Delta a_{0,1}^{(1)} = \sum_t (e_0^t \cdot a_{1,0}^{(2)} + e_1^t \cdot a_{1,1}^{(2)}) \alpha(O_1^{(2)t}) \cdot X_0^{(1)t}$$

$$\Delta a_{0,2}^{(1)} = \sum_t (e_0^t \cdot a_{2,0}^{(2)} + e_1^t \cdot a_{2,1}^{(2)}) \alpha(O_2^{(2)t}) \cdot X_0^{(1)t}$$

$$\Delta a_{0,3}^{(1)} = \sum_t (e_0 \cdot a_{3,0}^{(2)} + e_1^t \cdot a_{3,1}^{(2)}) \alpha(O_3^{(2)t}) \cdot X_0^{(1)t}$$

$$\Delta a_{1,1}^{(1)} = \sum_t (e_0^t \cdot a_{1,0}^{(2)} + e_1^t \cdot a_{1,1}^{(2)}) \alpha(O_1^{(2)t}) \cdot X_1^{(1)t}$$

$$\Delta a_{1,2}^{(1)} = \sum_t (e_0^t \cdot a_{2,0}^{(2)} + e_1^t \cdot a_{2,1}^{(2)}) \alpha(O_2^{(2)t}) \cdot X_1^{(1)t}$$

$$\Delta a^{(1)} = \begin{bmatrix} \Delta a_{0,1}^{(1)} & \Delta a_{0,2}^{(1)} & \Delta a_{0,3}^{(1)} & \Delta a_{0,4}^{(1)} \\ \Delta a_{1,1}^{(1)} & a_{1,2}^{(1)} & \Delta a_{1,3}^{(1)} & \Delta a_{1,4}^{(1)} \\ \Delta a_{2,1}^{(1)} & \Delta a_{2,2}^{(1)} & \Delta a_{2,3}^{(1)} & \Delta a_{2,4}^{(1)} \\ \vdots & & & \end{bmatrix}$$

$$\sum_t (e_0^t \cdot a_{1,0}^{(2)} + e_1^t \cdot a_{1,1}^{(2)}) \alpha(O_1^{(2)t}) \cdot X_0^{(1)t} \leq (e_0^t \cdot a_{2,0}^{(2)} + e_1^t \cdot a_{2,1}^{(2)}) \alpha(O_2^{(2)t}) \cdot X_0^{(1)t}$$

$$\sum_t (e_0^t \cdot a_{1,0}^{(2)} + e_1^t \cdot a_{1,1}^{(2)}) \alpha(O_1^{(2)t}) \cdot X_1^{(1)t} \leq (e_0^t \cdot a_{2,0}^{(2)} + e_1^t \cdot a_{2,1}^{(2)}) \alpha(O_2^{(2)t}) \cdot X_1^{(1)t}$$

$$\left. \begin{array}{c} e^t \cdot a_1^{(2)T} \alpha(O_1^{(2)t}) & e^t \cdot a_2^{(2)T} \alpha(O_2^{(2)t}) & e^t \cdot a_3^{(2)T} \alpha(O_3^{(2)t}) \\ e^t \cdot a_1^{(2)T} \alpha(O_1^{(2)t}) & e^t \cdot a_2^{(2)T} \alpha(O_2^{(2)t}) & e^t \cdot a_3^{(2)T} \alpha(O_3^{(2)t}) \\ \vdots & & \end{array} \right\}$$

$(e^t \cdot a^{(2)T}) \odot \alpha(O^{(2)t}) \cdot X_j^{(1)t}$

$$\Delta a_{j,k}^{(1)} = \left[ (e^t \cdot a^{(2)T}) \odot \alpha(O^{(2)t}) \cdot X_0^{(1)t} \right] \rightarrow \sum_t \begin{bmatrix} X_0^{(1)t} \\ X_1^{(1)t} \\ X_2^{(1)t} \\ X_3^{(1)t} \\ \vdots \end{bmatrix} \left[ (e^t \cdot a^{(2)T}) \odot \alpha(O^{(2)t}) \right]$$

$$\Delta a^{(1)} = \sum_t X^{(1)t} \cdot ((e^t \cdot a^{(2)T}) \odot \alpha(O^{(2)t}))$$

$$\begin{aligned}
 \Delta a_{i,j}^{(0)} &= \frac{\partial}{\partial a_{i,j}^{(0)}} \frac{\partial L}{\partial a_{i,j}^{(0)}} = \frac{\partial (\sum_{k=1}^m a_{k,m}^{(2)} X_k^{(2)t})}{\partial a_{i,j}^{(0)}} = \sum_{k=1}^m a_{k,m}^{(2)} \cdot \frac{\partial X_k^{(2)t}}{\partial a_{i,j}^{(0)}} = \frac{\partial X_k^{(2)t}}{\partial a_{i,j}^{(0)}} = \frac{\partial (\sigma(O_k^{(2)t}))}{\partial a_{i,j}^{(0)}} \\
 \frac{\partial \sigma(O_k^{(2)t})}{\partial a_{i,j}^{(0)}} &= \sigma(O_k^{(2)t})(1 - \sigma(O_k^{(2)t})) \frac{\partial O_k^{(2)t}}{\partial a_{i,j}^{(0)}} \frac{\partial O_k^{(2)t}}{\partial a_{i,j}^{(0)}} = \frac{\partial \left( \sum_{j=1}^n a_{j,k}^{(1)} \cdot X_j^{(1)t} \right)}{\partial a_{i,j}^{(0)}} \\
 &= a_{j,k}^{(1)} \cdot \frac{\partial X_j^{(1)t}}{\partial a_{i,j}^{(0)}} = \sigma(O_j^{(1)t})(1 - \sigma(O_j^{(1)t})) \cdot \frac{\partial O_j^{(1)t}}{\partial a_{i,j}^{(0)}} \frac{\partial O_j^{(1)t}}{\partial a_{i,j}^{(0)}} = \frac{\partial \left( \sum_{i=0}^l a_{i,j}^{(0)} \cdot X_i^{(0)t} \right)}{\partial a_{i,j}^{(0)}} \\
 &\Rightarrow \boxed{\Delta a_{i,j}^{(0)} = \gamma \sum_m \sum_{k=1}^m e_m^t \cdot \sum_{k=1}^m a_{k,m}^{(2)} \cdot \alpha(O_k^{(2)t}) \cdot a_{j,k}^{(1)} \cdot \alpha(O_j^{(1)t}) \cdot X_i^{(0)t}}
 \end{aligned}$$

vectorize in  $\Delta a_{i,j}^{(0)} = \gamma \sum_t [e_0^t \ e_1^t] \sum_{k=1}^m [a_{1,k}^{(2)} \ a_{2,k}^{(2)}] \cdot \cancel{[a_{3,k}^{(2)} \ a_{4,k}^{(2)}]} \cdot \cancel{\alpha(O_k^{(2)t})} \cdot a_{j,k}^{(1)} \cdot \alpha(O_j^{(1)t}) \cdot X_i^{(0)t}$

$$\Delta a_{i,j}^{(0)} = \gamma \sum_t [e_0^t \ e_1^t] \left[ [a_{1,0}^{(2)} \ a_{1,1}^{(2)}] \cdot \alpha(O_1^{(2)t}) \cdot a_{j,1} + [a_{2,0}^{(2)} \ a_{2,1}^{(2)}] \cdot \alpha(O_2^{(2)t}) \cdot a_{j,2} \right]$$

excluded +  $[a_{3,0}^{(2)} \ a_{3,1}^{(2)}] \cdot \alpha(O_3^{(2)t}) \cdot a_{j,3} + [a_{4,0}^{(2)} \ a_{4,1}^{(2)}] \cdot \alpha(O_4^{(2)t}) \cdot a_{j,4} + [a_{5,0}^{(2)} \ a_{5,1}^{(2)}] \cdot \alpha(O_5^{(2)t}) \cdot a_{j,5}$

$$a^{(2)} = \begin{bmatrix} a_{0,0} & a_{0,1} \\ a_{1,0} & a_{1,1} \\ a_{2,0} & a_{2,1} \\ a_{3,0} & a_{3,1} \\ a_{4,0} & a_{4,1} \\ a_{5,0} & a_{5,1} \end{bmatrix} \quad a^{(1)} = \begin{bmatrix} a_{1,1} & a_{1,2} & a_{1,3} & a_{1,4} & a_{1,5} \\ a_{2,1} & a_{2,2} & a_{2,3} & a_{2,4} & a_{2,5} \\ a_{3,1} & a_{3,2} & a_{3,3} & a_{3,4} & a_{3,5} \\ a_{4,1} & a_{4,2} & a_{4,3} & a_{4,4} & a_{4,5} \end{bmatrix}$$

$$\Delta a_{i,j}^{(0)} = \cancel{\sum_t [e_0^t \ e_1^t] \left[ a_{1,0}^{(2)} \ a_{1,1}^{(2)} \ a_{2,0}^{(2)} \ a_{2,1}^{(2)} \ a_{3,0}^{(2)} \ a_{3,1}^{(2)} \ a_{4,0}^{(2)} \ a_{4,1}^{(2)} \ a_{5,0}^{(2)} \ a_{5,1}^{(2)} \right] \circ \alpha(O^{(2)}) \cdot \begin{bmatrix} a_{j,1}^{(1)} \\ a_{j,2}^{(1)} \\ a_{j,3}^{(1)} \\ a_{j,4}^{(1)} \\ a_{j,5}^{(1)} \end{bmatrix} \cdot \alpha(O_j^{(1)t}) \cdot X_i^{(0)t}}$$

rebase j

$$\Delta a_{i,j}^{(0)} = \sum_t [e_0^t \ e_1^t] \left[ \begin{bmatrix} a_{1,0}^{(2)} & a_{2,0}^{(2)} & a_{3,0}^{(2)} & a_{4,0}^{(2)} & a_{5,0}^{(2)} \\ a_{1,1}^{(2)} & a_{2,1}^{(2)} & a_{3,1}^{(2)} & a_{4,1}^{(2)} & a_{5,1}^{(2)} \end{bmatrix} \circ \alpha(O^{(2)}) \cdot \begin{bmatrix} a_{1,1}^{(1)} & a_{2,1}^{(1)} & a_{3,1}^{(1)} & a_{4,1}^{(1)} & a_{5,1}^{(1)} \\ a_{1,2}^{(1)} & a_{2,2}^{(1)} & a_{3,2}^{(1)} & a_{4,2}^{(1)} & a_{5,2}^{(1)} \\ a_{1,3}^{(1)} & a_{2,3}^{(1)} & a_{3,3}^{(1)} & a_{4,3}^{(1)} & a_{5,3}^{(1)} \\ a_{1,4}^{(1)} & a_{2,4}^{(1)} & a_{3,4}^{(1)} & a_{4,4}^{(1)} & a_{5,4}^{(1)} \\ a_{1,5}^{(1)} & a_{2,5}^{(1)} & a_{3,5}^{(1)} & a_{4,5}^{(1)} & a_{5,5}^{(1)} \end{bmatrix} \cdot \alpha(O_1^{(1)t}) \cdot \alpha(O_2^{(1)t}) \cdot \dots \cdot \alpha(O_n^{(1)t}) \cdot X_i^{(0)t}$$

$$\Delta a_i^{(0)} = \sum_t \text{err. } a^{(2)\top} \circ \alpha(O^{(2)t}) \cdot a^{(1)\top} \circ \alpha(O_1^{(1)t}) \cdot X_i^{(0)t}$$

$$\Delta a^{(0)} = \sum_t \text{err. } a^{(2)\top} \circ \alpha(O^{(2)t}) \cdot a^{(1)\top} \circ \alpha(O^{(1)t}) \cdot \begin{bmatrix} X_0^{(0)t} \\ X_1^{(0)t} \\ X_2^{(0)t} \\ X_3^{(0)t} \end{bmatrix}$$

$$\boxed{\Delta a^{(0)} = \sum_t X^{(0)t\top} \cdot ((\text{err. } a^{(2)\top} \circ \alpha(O^{(2)t})) \cdot a^{(1)\top} \circ \alpha(O^{(1)t}))}$$

$$\Delta \mathbf{q}^{(2)} = \eta \cdot \mathbf{X}^{(2)T} \cdot \text{Err}^t \quad \text{-- only for one sample}$$

$$\Delta \mathbf{q}^{(1)} = \eta \cdot \mathbf{X}^{(1)T} \cdot ((\text{Err}^t \cdot \mathbf{q}^{(2)T}) \odot \alpha(\mathbf{O}^{(2)t}))$$

$$\Delta \mathbf{q}^{(0)} = \eta \sum_t \mathbf{X}^{(0)T} \cdot (((\text{Err} \cdot \mathbf{q}^{(2)T}) \odot \alpha(\mathbf{O}^{(2)t})) \cdot \mathbf{d}^{(1)T}) \odot \alpha(\mathbf{O}^{(1)t})$$

$\text{Err}^t \rightarrow$  Output layer error<sup>t</sup>

$$\text{Layer 2 error} \rightarrow ((\text{Output layer error}) \cdot \mathbf{q}^{(2)T}) \odot \alpha(\mathbf{O}^{(2)t})$$

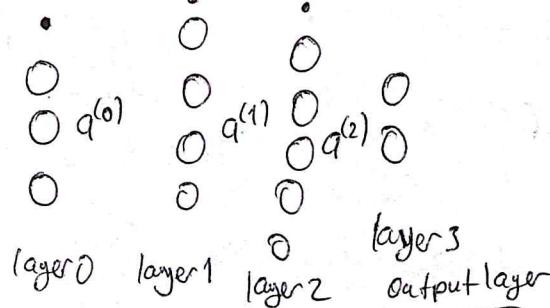
$$\text{Layer 1 error} \rightarrow ((\text{Layer 2 error}) \cdot \mathbf{q}^{(1)T}) \odot \alpha(\mathbf{O}^{(1)t})$$

$d^{(2)}_0 \rightarrow$  is excluded

$\downarrow$

$a^{(1)}_0 \rightarrow$  is excluded

$$\boxed{\text{Err}^{(l)t} = (\text{Err}^{(l+1)T} \cdot \mathbf{q}^{(l)T}) \odot \alpha(\mathbf{O}^{(l)t})}$$



$$\boxed{\Delta \mathbf{q}^{(l)} = \eta \cdot \mathbf{X}^{(l)T} \cdot \text{Err}^{(l+1)T}}$$

To deal with more than one sample we need to construct

a matrix  $\mathbf{X}$ .

$$\mathbf{X} = \begin{bmatrix} +1 & X_1^0 & X_2^0 & X_3^0 \\ +1 & X_1^1 & X_2^1 & X_3^1 \\ +1 & X_1^2 & X_2^2 & X_3^2 \end{bmatrix}$$

$$\mathbf{X} \cdot \mathbf{q}^{(0)} = \begin{bmatrix} +1 & O_1^0 & O_2^0 & O_3^0 & O_u^0 \\ +1 & O_1^1 & O_2^1 & O_3^1 & O_u^1 \\ +1 & O_1^2 & O_2^2 & O_3^2 & O_u^2 \end{bmatrix} \xrightarrow{\text{Layer added}} \begin{bmatrix} +1 & O_1^0 & O_2^0 & O_3^0 & O_4^0 & O_5^0 \\ +1 & O_1^1 & O_2^1 & O_3^1 & O_4^1 & O_5^1 \\ +1 & O_1^2 & O_2^2 & O_3^2 & O_4^2 & O_5^2 \end{bmatrix} \xrightarrow{\text{Layer 1 output}}$$

$$\boxed{(\text{Err} = r - o) \rightarrow \begin{bmatrix} e_0^0 & e_1^0 \\ e_0^1 & e_1^1 \\ e_0^2 & e_1^2 \end{bmatrix}}$$

$$\Delta \mathbf{q}^{(2)} = \eta \cdot \mathbf{X}^{(2)T} \cdot \text{Err}^t \quad \text{Err}^t \rightarrow \text{for a given sample}$$

$$\begin{bmatrix} +1 & X_1^{(2)0} & X_2^{(2)0} & X_3^{(2)0} & X_4^{(2)0} & X_5^{(2)0} \\ +1 & X_1^{(2)1} & X_2^{(2)1} & X_3^{(2)1} & X_4^{(2)1} & X_5^{(2)1} \\ +1 & X_1^{(2)2} & X_2^{(2)2} & X_3^{(2)2} & X_4^{(2)2} & X_5^{(2)2} \end{bmatrix} \rightarrow \mathbf{X}^{(2)} \text{ for } 3 \text{ samples}$$

$$\begin{aligned} L1\text{-Output} \cdot \mathbf{q}^{(1)} &= \begin{bmatrix} +1 & O_1^{(2)0} & O_2^{(2)0} & O_3^{(2)0} & O_4^{(2)0} & O_5^{(2)0} \\ +1 & O_1^{(2)1} & O_2^{(2)1} & O_3^{(2)1} & O_4^{(2)1} & O_5^{(2)1} \\ +1 & O_1^{(2)2} & O_2^{(2)2} & O_3^{(2)2} & O_4^{(2)2} & O_5^{(2)2} \end{bmatrix} \\ &\xleftarrow{\text{Layer 2 output}} \end{aligned}$$

$$\begin{aligned} \text{layer 2 Output} \cdot \mathbf{q}^{(2)} &= \begin{bmatrix} O_0^{(2)0} & O_1^{(2)0} \\ O_0^{(2)1} & O_1^{(2)1} \\ O_0^{(2)2} & O_1^{(2)2} \end{bmatrix} \\ &\xrightarrow{\text{Output layer}} \end{aligned}$$

$$\begin{aligned} \mathbf{q}^{(2)} &= \begin{bmatrix} +1 \\ X_1^{(2)0} \\ X_2^{(2)0} \\ X_3^{(2)0} \\ X_4^{(2)0} \\ X_5^{(2)0} \end{bmatrix} \begin{bmatrix} e_0^0 & e_1^0 \end{bmatrix} \begin{bmatrix} +1 \\ X_1^{(2)1} \\ X_2^{(2)1} \\ X_3^{(2)1} \\ X_4^{(2)1} \\ X_5^{(2)1} \end{bmatrix} \begin{bmatrix} e_0^1 & e_1^1 \end{bmatrix} \begin{bmatrix} +1 \\ X_1^{(2)2} \\ X_2^{(2)2} \\ X_3^{(2)2} \\ X_4^{(2)2} \\ X_5^{(2)2} \end{bmatrix} \begin{bmatrix} e_0^2 & e_1^2 \end{bmatrix} \\ &+ \begin{bmatrix} X_1^{(2)0} \\ X_2^{(2)0} \\ X_3^{(2)0} \\ X_4^{(2)0} \\ X_5^{(2)0} \end{bmatrix} \begin{bmatrix} +1 \\ X_1^{(2)1} \\ X_2^{(2)1} \\ X_3^{(2)1} \\ X_4^{(2)1} \\ X_5^{(2)1} \end{bmatrix} \begin{bmatrix} e_0^0 & e_1^0 \end{bmatrix} \begin{bmatrix} +1 \\ X_1^{(2)2} \\ X_2^{(2)2} \\ X_3^{(2)2} \\ X_4^{(2)2} \\ X_5^{(2)2} \end{bmatrix} \begin{bmatrix} e_0^1 & e_1^1 \end{bmatrix} \\ &+ \begin{bmatrix} X_1^{(2)1} \\ X_2^{(2)1} \\ X_3^{(2)1} \\ X_4^{(2)1} \\ X_5^{(2)1} \end{bmatrix} \begin{bmatrix} +1 \\ X_1^{(2)2} \\ X_2^{(2)2} \\ X_3^{(2)2} \\ X_4^{(2)2} \\ X_5^{(2)2} \end{bmatrix} \begin{bmatrix} e_0^0 & e_1^0 \end{bmatrix} \end{aligned}$$

$$\begin{aligned} L1\text{-Output} &= \sigma(L1\text{-Output}) \\ L2\text{-Output} &= \sigma(L2\text{-Output}) \end{aligned}$$

$$r = \begin{bmatrix} r_0^0 & r_1^0 \\ r_0^1 & r_1^1 \\ r_0^2 & r_1^2 \end{bmatrix} \quad \text{Output layer} = \text{softmax}(\text{Output layer}).$$

$$\Delta a^{(2)} = \begin{bmatrix} e_0^0 & e_1^0 \\ e_0^1 & e_1^1 \\ e_0^2 & e_1^2 \\ e_0^3 & e_1^3 \\ e_0^4 & e_1^4 \\ e_0^5 & e_1^5 \end{bmatrix} + \begin{bmatrix} e_0^0 & e_1^0 \\ e_0^1 & e_1^1 \\ e_0^2 & e_1^2 \\ e_0^3 & e_1^3 \\ e_0^4 & e_1^4 \\ e_0^5 & e_1^5 \end{bmatrix} + \begin{bmatrix} e_0^0 & e_1^0 \\ e_0^1 & e_1^1 \\ e_0^2 & e_1^2 \\ e_0^3 & e_1^3 \\ e_0^4 & e_1^4 \\ e_0^5 & e_1^5 \end{bmatrix}$$

$$\Delta a^{(2)} = \sum_t X_t^{(2)T} \cdot err^t$$

→ instead of multiplying individual input and one error we can combine them to matrix

$$\begin{bmatrix} 1 & 1 & 1 \\ X_1^{(2)0} & X_1^{(2)1} & X_1^{(2)2} \\ X_2^{(2)0} & X_2^{(2)1} & X_2^{(2)2} \\ X_3^{(2)0} & X_3^{(2)1} & X_3^{(2)2} \\ X_4^{(2)0} & X_4^{(2)1} & X_4^{(2)2} \\ X_5^{(2)0} & X_5^{(2)1} & X_5^{(2)2} \end{bmatrix} \begin{bmatrix} e_0^0 & e_1^0 \\ e_0^1 & e_1^1 \\ e_0^2 & e_1^2 \end{bmatrix} = \begin{bmatrix} (e_0^0 + e_0^1 + e_0^2) & (e_1^0 + e_1^1 + e_1^2) \\ (X_1^{(2)0} \cdot e_0^0 + X_1^{(2)1} \cdot e_0^1 + X_1^{(2)2} \cdot e_0^2) & (X_1^{(2)0} \cdot e_1^0 + X_1^{(2)1} \cdot e_1^1 + X_1^{(2)2} \cdot e_1^2) \\ (X_2^{(2)0} \cdot e_0^0 + X_2^{(2)1} \cdot e_0^1 + X_2^{(2)2} \cdot e_0^2) & (X_2^{(2)0} \cdot e_1^0 + X_2^{(2)1} \cdot e_1^1 + X_2^{(2)2} \cdot e_1^2) \\ \vdots & \vdots \\ \vdots & \vdots \end{bmatrix}$$

$$\Delta a^{(1)} = \sum_t X^{(1)T} t \cdot ((err^t \cdot a^{(2)T}) \otimes (\alpha^{(2)t}))$$

$$err^t = \begin{bmatrix} e_0^t & e_1^t \end{bmatrix} \quad a^{(2)} = \begin{bmatrix} a_1^{(2)} \\ a_2^{(2)} \\ a_3^{(2)} \\ a_4^{(2)} \\ a_5^{(2)} \end{bmatrix} \quad a_i^{(2)} = \begin{bmatrix} a_{i,1}^{(2)} & a_{i,2}^{(2)} \end{bmatrix}$$

$$\begin{bmatrix} -e_0^t \\ -e_1^t \end{bmatrix} \cdot \begin{bmatrix} a_1^{(2)} & a_2^{(2)} & a_3^{(2)} & a_4^{(2)} & a_5^{(2)} \end{bmatrix} \otimes \begin{bmatrix} \alpha^{(2)0} \\ \alpha^{(2)1} \\ \alpha^{(2)2} \end{bmatrix}$$

$$\begin{bmatrix} (e_0^t a_1^{(2)}) (e_0^t a_2^{(2)}) (e_0^t a_3^{(2)}) (e_0^t a_4^{(2)}) (e_0^t a_5^{(2)}) \\ e_1^t a_1^{(2)} & e_1^t a_2^{(2)} & e_1^t a_3^{(2)} & e_1^t a_4^{(2)} & e_1^t a_5^{(2)} \\ e^2 a_1^{(2)} & e^2 a_2^{(2)} & e^2 a_3^{(2)} & e^2 a_4^{(2)} & e^2 a_5^{(2)} \end{bmatrix} \otimes \begin{bmatrix} \alpha^{(2)0} \\ \alpha^{(2)1} \\ \alpha^{(2)2} \end{bmatrix}$$

$$X^{(1)T} = \begin{bmatrix} 1 & 1 & 1 \\ X_1^{(1)0} & X_1^{(1)1} & X_1^{(1)2} \\ 1 & 1 & 1 \end{bmatrix}$$

$$\Delta a^{(1)} = X^{(1)T} \cdot$$

$$\beta_x^t = e_0^t a_x^{(2)} \alpha^{(2)0}$$

$$\Delta a^{(1)} = \begin{bmatrix} (\beta_1^0 \cdot 1 + \beta_1^1 \cdot 1 + \beta_1^2) \\ (\beta_1^0 \cdot X_1^{(1)0} + \beta_1^1 \cdot X_1^{(1)1} + \beta_1^2 \cdot X_1^{(1)2}) \\ (\beta_1^0 \cdot X_2^{(1)0} + \beta_1^1 \cdot X_2^{(1)1} + \beta_1^2 \cdot X_2^{(1)2}) \\ (\beta_1^0 \cdot X_3^{(1)0} + \beta_1^1 \cdot X_3^{(1)1} + \beta_1^2 \cdot X_3^{(1)2}) \end{bmatrix}$$

$$\begin{bmatrix} (\beta_2^0 \cdot a_1^{(2)0}) e_0^t a_2^{(2)} \alpha^{(2)0} & e_0^t a_2^{(2)} \alpha^{(2)0} e_0^t a_3^{(2)} \alpha^{(2)0} (e_0^t a_4^{(2)} \alpha^{(2)0}) \\ e_1^t a_1^{(2)} \alpha^{(2)1} e_1^t a_2^{(2)} \alpha^{(2)1} e_0^t a_3^{(2)} \alpha^{(2)1} (e_0^t a_4^{(2)} \alpha^{(2)1}) \\ e^2 a_1^{(2)} \alpha^{(2)2} e_2^t a_2^{(2)} \alpha^{(2)2} e_0^t a_3^{(2)} \alpha^{(2)2} (e_0^t a_4^{(2)} \alpha^{(2)2}) \end{bmatrix}$$

↓ layer 2 error

$$\begin{bmatrix} (\beta_2^0 + \beta_2^1 + \beta_2^2) (\beta_3^0 + \beta_3^1 + \beta_3^2) (\beta_4^0 + \beta_4^1 + \beta_4^2) (\beta_5^0 + \beta_5^1 + \beta_5^2) \\ (\beta_2^0 \cdot X_1^{(1)0} + \beta_2^1 \cdot X_1^{(1)1} + \beta_2^2 \cdot X_1^{(1)2}) \\ \beta_2^0 \cdot X_2^{(1)0} + \beta_2^1 \cdot X_2^{(1)1} + \beta_2^2 \cdot X_2^{(1)2} \end{bmatrix}$$

$$\begin{aligned} \Delta a_{j,k}^{(1)} &= \eta \sum_t \sum_m e_m^t \cdot q_{k,m}^{(2)} \cdot \alpha(O_k^{(2)t}) \cdot X_j^{(1)t} \\ &= \eta \sum_t (\underbrace{e^t \cdot a_k^{(2)}}_{\substack{\text{defined} \\ \text{previously}}} \cdot \alpha(O_k^{(2)t}) \cdot X_j^{(1)t}) \end{aligned}$$

$$\Delta a_{j,k}^{(1)} = \eta \cdot \begin{bmatrix} \sum_t \beta_0^t \cdot X_0^{(1)t} & \sum_t \beta_1^t \cdot X_0^{(1)t} & \sum_t \beta_2^t \cdot X_0^{(1)t} & \sum_t \beta_3^t \cdot X_0^{(1)t} & \sum_t \beta_4^t \cdot X_0^{(1)t} \\ \sum_t \beta_1^t \cdot X_1^{(1)t} & \sum_t \beta_2^t \cdot X_1^{(1)t} & \sum_t \beta_3^t \cdot X_1^{(1)t} & \sum_t \beta_4^t \cdot X_1^{(1)t} & \dots \\ \sum_t \beta_2^t \cdot X_2^{(1)t} & \sum_t \beta_3^t \cdot X_2^{(1)t} & \sum_t \beta_4^t \cdot X_2^{(1)t} & \dots & \dots \\ \sum_t \beta_3^t \cdot X_3^{(1)t} & \dots & \dots & \dots & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

Same thing represented on previous page.

$$\Delta a_{i,j}^{(0)} = \eta \sum_t \sum_m e_m^t \cdot \sum_{k=1}^m q_{k,m}^{(2)} \underbrace{a_k^{(2)t} \cdot \alpha(O_k^{(2)t}) \cdot a_{j,k}^{(1)} \cdot \alpha(O_j^{(1)t}) \cdot X_i^{(0)t}}_{\substack{\text{layer 2 error} \\ \text{fwd will produce}}}$$

$$\Delta a^{(0)} = \sum_t X^{(0)T} \left( \frac{\text{err. } a^{(2)T} \alpha(O^{(2)T}) \cdot a^{(1)T}}{\alpha(O^{(1)T})} \right)$$

If we represent layer 2 error  $\begin{bmatrix} e^{(0)} \\ e^{(1)} \\ e^{(2)} \end{bmatrix}$

$$\Delta a^{(0)} = \sum_t X^{(0)T} \cdot \left[ \frac{\text{err. of prev layer} \cdot a^{(1)T}}{\alpha(O^{(1)T})} \right] \alpha(O^{(1)T})$$

$$X^{(0)T} = \begin{bmatrix} 1 & 1 & 1 \\ X^{(0)0} & X^{(0)1} & X^{(0)2} \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 1 \\ a_1^{(1)} & a_2^{(1)} & a_3^{(1)} & a_4^{(1)} \end{bmatrix} \odot \begin{bmatrix} \alpha(O^{(1)0}) \\ \alpha(O^{(1)1}) \\ \alpha(O^{(1)2}) \end{bmatrix}$$

$$\begin{aligned} X^{(0)T} \cdot \alpha(O^{(1)T}) &= \begin{bmatrix} e^{(0)a_1^{(1)}} \alpha(O_1^{(1)0}) & e^{(0)a_2^{(1)}} \alpha(O_2^{(1)0}) & e^{(0)a_3^{(1)}} \alpha(O_3^{(1)0}) & e^{(0)a_4^{(1)}} \alpha(O_4^{(1)0}) \\ e^{(1)a_1^{(1)}} \alpha(O_1^{(1)1}) & e^{(1)a_2^{(1)}} \alpha(O_2^{(1)1}) & e^{(1)a_3^{(1)}} \alpha(O_3^{(1)1}) & e^{(1)a_4^{(1)}} \alpha(O_4^{(1)1}) \\ e^{(2)a_1^{(1)}} \alpha(O_1^{(1)2}) & e^{(2)a_2^{(1)}} \alpha(O_2^{(1)2}) & e^{(2)a_3^{(1)}} \alpha(O_3^{(1)2}) & e^{(2)a_4^{(1)}} \alpha(O_4^{(1)2}) \end{bmatrix} \\ \text{Again } \bar{\beta}_x^t &= e^{(t)} a_x^{(1)} \alpha(O_x^{(1)0}) \quad (\bar{\beta}_0 + \bar{\beta}_1 + \bar{\beta}_2) \quad (\bar{\beta}_3 + \bar{\beta}_4 + \bar{\beta}_5) \quad (\bar{\beta}_6 + \bar{\beta}_7 + \bar{\beta}_8) \end{aligned}$$

$$= \begin{bmatrix} \bar{\beta}_0 \cdot X_1^{(0)0} + \bar{\beta}_1 \cdot X_1^{(0)1} + \bar{\beta}_2 \cdot X_1^{(0)2} \\ \bar{\beta}_3 \cdot X_2^{(0)0} + \bar{\beta}_4 \cdot X_2^{(0)1} + \bar{\beta}_5 \cdot X_2^{(0)2} \\ \bar{\beta}_6 \cdot X_3^{(0)0} + \bar{\beta}_7 \cdot X_3^{(0)1} + \bar{\beta}_8 \cdot X_3^{(0)2} \end{bmatrix}$$

$$\Delta a_{i,j}^{(0)} = \eta \sum_t e_k^t a_{j,k} \alpha(O_j^{(1)t}) \cdot X_i^{(0)t}$$

$$\Delta a_{i,j}^{(0)} = \eta \sum_t \frac{(e_k^t \cdot a_j^{(1)}) \alpha(O_j^{(1)t}) \cdot X_i^{(0)t}}{\bar{\beta}_j^t}$$

$$\Delta a_i^{(0)} = \eta \begin{bmatrix} \sum_t \bar{\beta}_1^t \cdot X_0^{(0)t} & \sum_t \bar{\beta}_2^t \cdot X_0^{(0)t} & \sum_t \bar{\beta}_3^t \cdot X_0^{(0)t} & \sum_t \bar{\beta}_4^t \cdot X_0^{(0)t} \\ \sum_t \bar{\beta}_1^t \cdot X_1^{(0)t} & \sum_t \bar{\beta}_2^t \cdot X_1^{(0)t} & \sum_t \bar{\beta}_3^t \cdot X_1^{(0)t} & \sum_t \bar{\beta}_4^t \cdot X_1^{(0)t} \\ \sum_t \bar{\beta}_1^t \cdot X_2^{(0)t} & \sum_t \bar{\beta}_2^t \cdot X_2^{(0)t} & \sum_t \bar{\beta}_3^t \cdot X_2^{(0)t} & \sum_t \bar{\beta}_4^t \cdot X_2^{(0)t} \end{bmatrix}$$

Same as  
on prev. page

It can also be shown that the same procedure can be applied  
to inner layers. By subtracting matrix of inputs and output labels  
for all samples. By subtracting matrix of inputs and output labels  
Batch of samples could be used to update weight matrices  
simultaneously in one pass with matrix multiplications.

In previous we have shown that

$$\Delta a^{(2)} = \eta \cdot X^{(2)T} \cdot Err^t$$

$$\Delta a^{(1)} = \eta \cdot X^{(1)T} \cdot ((Err^t \cdot a^{(2)T}) \odot \alpha(O^{(2)t}))$$

$$\Delta a^{(0)} = \eta \cdot X^{(0)T} \cdot (((((Err^t \cdot a^{(2)T}) \odot \alpha(O^{(2)t})) \cdot a^{(1)T}) \odot \alpha(O^{(1)t}))$$

$\Rightarrow$  There are for one sample

so as we have shown previously we can stack samples in a

matrix

$$X = \begin{bmatrix} -x_0^0 \\ -x_0^1 \\ -x_0^2 \\ -x_0^3 \end{bmatrix}$$

$$\text{where } X^{(0)t} = \begin{bmatrix} X_0^{(0)t} & X_1^{(0)t} & X_2^{(0)t} & X_3^{(0)t} \end{bmatrix}$$

$$r = \begin{bmatrix} -r^0 \\ -r^1 \\ -r^2 \\ -r^3 \end{bmatrix} \quad \text{labels}$$

$$r^t = [r_0^t \ r_1^t]$$

so. for batch

$$\boxed{\Delta a^{(2)} = \eta \cdot X^{(2)T} \cdot Err^t}$$

$$\boxed{\Delta a^{(1)} = \eta \cdot X^{(1)T} \cdot (((Err^t \cdot a^{(2)T}) \odot \alpha(O^{(2)t})), a^{(1)T}) \odot \alpha(O^{(1)t})}$$

$$\boxed{\Delta a^{(0)} = \eta \cdot X^{(0)T} \cdot ((Err^t \cdot a^{(2)T}) \odot \alpha(O^{(2)t}))}$$

$$\boxed{Err = O^{(3)} - r}$$

$$\text{Output error} = O^{(3)} - r$$

$$\text{Layer 2 error} = ((\text{Output Error}), a^{(2)T}) \odot \alpha(O^{(2)})$$

$$\text{Layer 1 error} = ((\text{Layer 2 error}), a^{(1)T}) \odot \alpha(O^{(1)})$$

$$\text{Error}^{(e)} = (Err^{(e+1)}, a^{(e)T}) \odot \alpha(O^{(e)})$$

$$\Delta a^{(2)} = \eta \cdot X^{(2)T} \cdot \text{Output Error}$$

$$\Delta a^{(1)} = \eta \cdot X^{(1)T} \cdot \text{Layer 2 error}$$

$$\Delta a^{(0)} = \eta \cdot X^{(0)T} \cdot \text{Layer 1 error}$$

$$\Delta a^{(e)} = \eta \cdot X^{(e)T} \cdot Err^{(e+1)}$$

→ operation is used for matrix multiplication, scalar and matrix multiplication, scalar-scalar multiplication

② → operation is used for matNx-matrix, vector-vector elementwise multiplication

$$X = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \quad Y = \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} \quad X \cdot Y = \begin{bmatrix} 19 & 22 \\ 43 & 50 \end{bmatrix} \quad X \odot Y = \begin{bmatrix} 1 \times 5 & 2 \times 6 \\ 3 \times 7 & 4 \times 8 \end{bmatrix}$$

$$\Delta a^{(2)} = \eta \cdot X^{(2)T} \cdot Err^T$$

↓  
scalar matrix multiplication

↓  
matrix multiplication

↓  
matrix multiplication → pairwise multiplication

↓  
pairwise multiplication

$$\Delta a^{(1)} = \eta \cdot X^{(1)T} \cdot ((Err \cdot a^{(2)T}) \odot \alpha(O^{(2)}))$$

- Parenthesis has the most high priority in calculation. First the terms inside the parenthesis are handled.

$$\Delta a^{(0)} = \eta \cdot X^{(0)T} \cdot (((((Err \cdot a^{(2)T}) \odot \alpha(O^{(2)})) \cdot a^{(1)T}) \odot \alpha(O^{(1)}))$$

↓  
scalar-matrix multiplication

↓  
matrix multiplication

↓  
pairwise multiplication

- sigmoid( $x$ ) =  $\sigma(x)$        $\sigma(x) = \frac{1}{1 + e^{-x}}$

$$(\sigma(x))^t = ((1 + e^{-x})^{-1})^t = -1 \cdot \frac{(1 + e^{-x})^t}{(1 + e^{-x})^2} = \frac{+1 \cdot (e^{-x} \cdot +1)}{(1 + e^{-x})^2}$$

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad \sigma(x) + \sigma(x) \cdot e^{-x} = 1 \quad e^{-x} = \frac{1 - \sigma(x)}{\sigma(x)}$$

$$\sigma(x)^t = \frac{e^{-x}}{(1 + e^{-x})^2} = \sigma(x)^2 \cdot e^{-x} = \sigma(x)^2 \cdot \frac{(1 - \sigma(x))}{\sigma(x)} = \sigma(x) \cdot (1 - \sigma(x))$$

$\sigma(x)^t = \sigma(x) \cdot (1 - \sigma(x))$

Classification  
cost function  
Derivation

Consider we have  $n$  number of class in the given data set.  
At the output of the network we use softmax to obtain the probability distribution of each class. For each class our network a probability  $P_i$  for ~~each class~~. We can think of this distribution as multinomial density.

$P_i = P(c_i | \theta)$  represents our network parameters (out network).

So we want our network to produce maximum likelihood estimation for each class in dataset. (Maximum Likelihood Estimation.)

$$P(c_1, c_2, c_3, \dots, c_t) = \prod_{i=1}^t P_i^{c_i}$$

probability of whole data set classes.

$$c_i = \begin{cases} \text{████████} & \text{if current label is } i \\ 1 & \text{otherwise} \end{cases}$$

→ because each class probability is independent.

$$P(c_1, c_2, c_3, \dots, c_t) = \prod_{t=1}^T \prod_{i=1}^{n_t} P(c_i) = \prod_{t=1}^T \prod_{i=1}^{n_t} P(k_i | \theta)$$

so we are looking for the optimum  $\theta$  value which maximizes the probability.  
We could take log of this expression because max value location will not change.

$$P(\text{All set}) = \log \prod_{t=1}^T \prod_{i=1}^{n_t} P(k_i | \theta)^{c_i} = \sum_{t=1}^T \sum_{i=1}^{n_t} c_i \log P(k_i | \theta)$$

This maximization problem but by ~~maximizing~~ flipping maximization problem to minimization problem so

$$E = j = - \sum_{t=1}^T \sum_{m=1}^{n_t} r_m \cdot \log(y_m^t) \quad \rightarrow \text{neural network output.}$$

$$\Delta \alpha_{k,m}^{(2)} = -\eta \frac{\partial E}{\partial \alpha_{k,m}^{(2)}} \quad \frac{\partial E}{\partial \alpha_{k,m}^{(2)}} = \left( - \sum_{t=1}^T \sum_{m=1}^{n_t} r_m \cdot \log(y_m^t) \right) \frac{\partial \log(y_m^t)}{\partial \alpha_{k,m}^{(2)}}$$

$$= \sum_{t=1}^T \sum_{m=1}^{n_t} -r_m \frac{\log(y_m^t)}{\partial \alpha_{k,m}^{(2)}} \quad \frac{\partial \log(y_m^t)}{\partial \alpha_{k,m}^{(2)}} = \frac{1}{y_m^t} \cdot \frac{\partial y_m^t}{\partial \alpha_{k,m}^{(2)}}$$

$$y_m^t = \frac{e^{O_m^{(3)}}}{\sum_{n=1}^N e^{O_n^{(3)}}}$$

$$[y_m = \text{softmax}(O_m^{(3)})]$$

$$\text{softmax}(i) = \frac{e^i}{\sum_n e^n}$$

$$\frac{d \text{softmax}(i)}{d i} = \frac{1}{\sum_n e^n} \cdot e^i - \frac{1}{(\sum_n e^n)^2} e^i \cdot e^i$$

$$= \frac{e^i \sum_n e^n - e^{2i}}{(\sum_n e^n)^2}$$

$$= e^i \cdot \left( \sum_n e^n - e^i \right)$$

$$= \frac{e^i}{(\sum_n e^n)} \cdot \left( \sum_n e^n - \frac{e^i}{\sum_n e^n} \right)$$

$$= \text{softmax}(i) \cdot \left( 1 - \text{softmax}(i) \right)$$

$$\sum_{t=1}^T \sum_{i=1}^{n_t} c_i \log(P(k_i | \theta))$$

$P(k_i | \theta)$  represents the probability of any class at that given sample data.  $c_i$  represents which class is the correct one.

Simply our network output a probability distribution

$$[0.1 \ 0.2 \ 0.3 \ 0.4] \quad \text{class A} \ \text{class B} \ \text{class C} \ \text{class D}$$

$c_j$  represents if given data

$$j \text{ class A} \quad c_j = [1 \ 0 \ 0 \ 0]$$

if class B

$$c_j = [0 \ 1 \ 0 \ 0]$$

if class C

$$c_j = [0 \ 0 \ 1 \ 0]$$

$$\text{if } D \quad c_j = [0 \ 0 \ 0 \ 1]$$

when correct class' probability is less this causes error to be high. for ex.

$$[0.9 \ 0.09 \ 0.009 \ 0.001] \quad \text{and correct class be } [0 \ 0 \ 0 \ 1]$$

$$c_i \cdot \log([0.9 \ 0.09 \ 0.009 \ 0.001])$$

$$[0 \ 0 \ 0 \ 1] \cdot \log([0.9 \ 0.09 \ 0.009 \ 0.001])$$

$$\log(0.001)$$

this is a negative value

but at front there is also one  $-1$ . so the error becomes positive and higher.

$$\frac{d \text{softmax}(i)}{d j} = -\frac{e^j \cdot e^i}{(\sum_n e^n)^2} = -\frac{e^j}{\sum_n e^n} \cdot \frac{e^i}{\sum_n e^n} = -\text{softmax}(i) \cdot \text{softmax}(j)$$

$$\frac{d \text{softmax}(i)}{d j} = \begin{cases} \text{softmax}(i) \cdot (1 - \text{softmax}(j)) & \text{if } i=j \\ -\text{softmax}(i) \cdot \text{softmax}(j) & \text{otherwise} \end{cases}$$

$$E = j = \sum_t \sum_m r_m t \cdot \log(y_m t)$$

$$\frac{\partial E}{\partial a_{k,n}^{(2)}} = ?$$

$$x_0^{(2)} = +1$$

$$y_m t = \text{softmax}(o_m^{(3)t})$$

$$\frac{\partial E}{\partial a_{k,n}^{(2)}} = -\sum_t \sum_m r_m t \frac{\partial \log(y_m t)}{\partial a_{k,n}^{(2)}}$$

$$\begin{aligned} \frac{\partial \log(y_m t)}{\partial a_{k,n}^{(2)}} &= \frac{\partial \log y_m t}{\partial y_m t} \cdot \frac{\partial y_m t}{\partial o_n^{(2)t}} \cdot \frac{\partial o_n^{(2)t}}{\partial a_{k,n}^{(2)}} \\ &= \frac{1}{y_m t} \cdot y_m t (1_{\{m=n\}} - y_n t) \cdot \frac{\partial o_n^{(2)t}}{\partial a_{k,n}^{(2)}} \end{aligned}$$

layer 2 layer 3

$k \quad n \in M$

layer 2 layer 3

$y_m t$  is function of  $o_n^{(2)t}$  at the same time.  $y_m t = \text{softmax}(o_m^{(2)t}) = \frac{e^{o_m^{(2)t}}}{\sum_k e^{o_k^{(2)t}}}$

$n$  represents the output layer.

$$\begin{aligned} \frac{\partial E}{\partial a_{k,n}^{(2)}} &= -\sum_t \sum_m r_m t \frac{1}{y_m t} \cdot y_m t (1_{\{m=n\}} - y_n t) \cdot \frac{\partial o_n^{(2)t}}{\partial a_{k,n}^{(2)}} \\ &= -\sum_t \sum_m r_m t (1_{\{m=n\}} - y_n t) \frac{\partial o_n^{(2)t}}{\partial a_{k,n}^{(2)}} \\ &= -\sum_t r_{m=n}^t (1 - y_n t) \frac{\partial o_n^{(2)t}}{\partial a_{k,n}^{(2)}} + \sum_{m \neq n} r_m t \cdot y_n t \cdot \frac{\partial o_n^{(2)t}}{\partial a_{k,n}^{(2)}} \\ &= -\sum_t r_{m=n}^t \frac{\partial o_n^{(2)t}}{\partial a_{k,n}^{(2)}} - r_{m=n}^t y_n t \frac{\partial o_n^{(2)t}}{\partial a_{k,n}^{(2)}} + -y_n t \sum_{m \neq n} r_m t \cdot \frac{\partial o_n^{(2)t}}{\partial a_{k,n}^{(2)}} \\ &= -\sum_t r_{m=n}^t \frac{\partial o_n^{(2)t}}{\partial a_{k,n}^{(2)}} - y_n t \cdot \frac{\partial o_n^{(2)t}}{\partial a_{k,n}^{(2)}} \cdot (+r_m^t + \sum_{m \neq n} r_m t) \rightarrow r^t = [0 \ 0 \ 0 \ 1] = \text{sum} = 1 \\ &\quad \boxed{-\sum_t (r_n^t - y_n t) \frac{\partial o_n^{(2)t}}{\partial a_{k,n}^{(2)}}} \quad \text{or } [0 \ 1 \ 0 \ 0] = \text{sum} = 1 \end{aligned}$$

$$\Delta a_{k,n}^{(2)} = -\eta \frac{\partial E}{\partial a_{k,n}^{(2)}} = \eta \sum_t (r_n^t - y_n t) \frac{\partial o_n^{(2)t}}{\partial a_{k,n}^{(2)}}$$

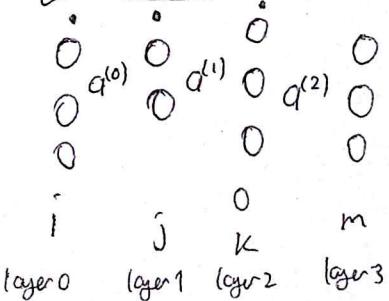
$$\text{In regression } E = j = \sum_t \sum_m (r_m t - o_m^{(3)t})^2 \rightarrow o_m^{(3)t} = \sum_k a_{k,m}^{(2)} x_k^{(2)t}$$

$$\Delta a_{k,n}^{(2)} = -\eta \frac{\partial E}{\partial a_{k,n}^{(2)}} \quad \frac{\partial E}{\partial a_{k,n}^{(2)}} = \sum_t \sum_m (r_m t - o_m^{(3)t}) \cdot 1 \cdot \frac{\partial o_m^{(3)t}}{\partial a_{k,n}^{(2)}}$$

$$\Delta a_{k,n}^{(2)} = \boxed{+\eta \sum_t \sum_m (r_m t - o_m^{(3)t}) + 1 \cdot \frac{\partial o_m^{(3)t}}{\partial a_{k,n}^{(2)}}}$$

in regression and classification loss functions' derivations of back propagation algorithm are same.

## Implementation Considerations



In forward pass.

$$\left[ \begin{array}{ccc} (+1) & x_0^{(0)} & x_1^{(0)} \\ (+1) & x_1^{(0)} & x_2^{(0)} \\ (+1) & x_2^{(0)} & x_3^{(0)} \end{array} \right] \cdot \left[ \begin{array}{cc} a_{0,1}^{(0)} & a_{0,2}^{(0)} \\ a_{1,1}^{(0)} & a_{1,2}^{(0)} \\ a_{2,1}^{(0)} & a_{2,2}^{(0)} \\ a_{3,1}^{(0)} & a_{3,2}^{(0)} \end{array} \right]$$

*Note: This part should be inserted and they correspond to bias values multiplier (+1)*

The result of this multiplication

For every eight input multiplication we are applying linear transformation

$$(+1) \left[ \begin{array}{cc} O_1^{(1)0} & O_2^{(1)0} \\ O_1^{(1)1} & O_2^{(1)1} \end{array} \right] \cdot \left[ \begin{array}{cccc} a_{0,1}^{(1)} & a_{0,2}^{(1)} & a_{0,3}^{(1)} & a_{0,4}^{(1)} \\ a_{1,1}^{(1)} & a_{1,2}^{(1)} & a_{1,3}^{(1)} & a_{1,4}^{(1)} \\ a_{2,1}^{(1)} & a_{2,2}^{(1)} & a_{2,3}^{(1)} & a_{2,4}^{(1)} \\ a_{3,1}^{(1)} & a_{3,2}^{(1)} & a_{3,3}^{(1)} & a_{3,4}^{(1)} \end{array} \right]$$

for this multiplication to be reformed this column vector should be added. And these are multiplier of biases (+1).

This results in

$$(+1) \left[ \begin{array}{cccc} O_1^{(2)0} & O_2^{(2)0} & O_3^{(2)0} & O_4^{(2)0} \\ O_1^{(2)1} & O_2^{(2)1} & O_3^{(2)1} & O_4^{(2)1} \end{array} \right]$$

To perform this operation we must add a column vector

$$\left[ \begin{array}{ccc} a_{0,0}^{(2)} & a_{0,1}^{(2)} & a_{0,2}^{(2)} \\ a_{1,0}^{(2)} & a_{1,1}^{(2)} & a_{1,2}^{(2)} \\ a_{2,0}^{(2)} & a_{2,1}^{(2)} & a_{2,2}^{(2)} \\ a_{3,0}^{(2)} & a_{3,1}^{(2)} & a_{3,2}^{(2)} \\ a_{4,0}^{(2)} & a_{4,1}^{(2)} & a_{4,2}^{(2)} \end{array} \right]$$

So to get rid of computation overhead of column insertion, we can separate each weight matrix's first row. Because they represent the bias of weights.

$$a^{(0)} = \left[ \begin{array}{cc} a_{0,1}^{(0)} & a_{0,2}^{(0)} \\ a_{1,1}^{(0)} & a_{1,2}^{(0)} \\ a_{2,1}^{(0)} & a_{2,2}^{(0)} \\ a_{3,1}^{(0)} & a_{3,2}^{(0)} \end{array} \right] \rightarrow \left[ \begin{array}{cc} a_{1,1}^{(0)} & a_{1,2}^{(0)} \\ a_{2,1}^{(0)} & a_{2,2}^{(0)} \\ a_{3,1}^{(0)} & a_{3,2}^{(0)} \end{array} \right] \text{ and } \left[ \begin{array}{c} a_{0,1}^{(0)} \\ a_{0,2}^{(0)} \end{array} \right]$$

new  $a^{(0)}$

$b^{(0)}$

$$X = \begin{bmatrix} x_0^{(0)} & x_1^{(0)} & x_2^{(0)} & x_3^{(0)} \end{bmatrix}$$

$$O_j^{(1)t} = \sum_{i=0}^3 a_{i,j}^{(0)} X_i^{(0)t}$$

↓ input form

$$O^{(1)} = X \cdot a^{(0)}$$

$$a^{(0)} = \begin{bmatrix} a_{0,1}^{(0)} & a_{0,2}^{(0)} \\ a_{1,1}^{(0)} & a_{1,2}^{(0)} \\ a_{2,1}^{(0)} & a_{2,2}^{(0)} \\ a_{3,1}^{(0)} & a_{3,2}^{(0)} \end{bmatrix}$$

$$X_j^{(1)t} = \sigma(O_j^{(1)t})$$

$$\sigma(O^{(1)}) = X^{(1)}$$

$$O_k^{(2)t} = \sum_{j=0}^3 a_{j,k}^{(1)} X_j^{(1)t}$$

$$O^{(2)} = X^{(1)} \cdot a^{(1)}$$

$$a^{(1)} = \begin{bmatrix} a_{0,1}^{(1)} & a_{0,2}^{(1)} & a_{0,3}^{(1)} & a_{0,4}^{(1)} \\ a_{1,1}^{(1)} & a_{1,2}^{(1)} & a_{1,3}^{(1)} & a_{1,4}^{(1)} \\ a_{2,1}^{(1)} & a_{2,2}^{(1)} & a_{2,3}^{(1)} & a_{2,4}^{(1)} \\ a_{3,1}^{(1)} & a_{3,2}^{(1)} & a_{3,3}^{(1)} & a_{3,4}^{(1)} \end{bmatrix}$$

$$X^{(2)} = \sigma(O^{(2)})$$

$$O_m^{(3)t} = \sum_{k=0}^4 a_{k,m}^{(2)} X_k^{(2)t}$$

$$\text{Output} = O^{(3)} = X^{(2)} \cdot a^{(2)} \quad \text{for regression}$$

$$a^{(2)} = \begin{bmatrix} a_{0,0}^{(2)} & a_{0,1}^{(2)} & a_{0,2}^{(2)} \\ a_{1,0}^{(2)} & a_{1,1}^{(2)} & a_{1,2}^{(2)} \\ a_{2,0}^{(2)} & a_{2,1}^{(2)} & a_{2,2}^{(2)} \\ a_{3,0}^{(2)} & a_{3,1}^{(2)} & a_{3,2}^{(2)} \\ a_{4,0}^{(2)} & a_{4,1}^{(2)} & a_{4,2}^{(2)} \end{bmatrix}$$

In classification

$$\text{Output} = \text{softmax}(O^{(3)})$$

$$q^{(1)} = \begin{bmatrix} q_{0,1}^{(1)} & q_{0,2}^{(1)} & q_{0,3}^{(1)} & q_{0,4}^{(1)} & q_{0,5}^{(1)} \\ q_{1,1}^{(1)} & q_{1,2}^{(1)} & q_{1,3}^{(1)} & q_{1,4}^{(1)} & q_{1,5}^{(1)} \\ q_{2,1}^{(1)} & q_{2,2}^{(1)} & q_{2,3}^{(1)} & q_{2,4}^{(1)} & q_{2,5}^{(1)} \end{bmatrix} \rightarrow \begin{bmatrix} q_{1,1}^{(1)} & q_{1,2}^{(1)} & q_{1,3}^{(1)} & q_{1,4}^{(1)} & q_{1,5}^{(1)} \\ q_{2,1}^{(1)} & q_{2,2}^{(1)} & q_{2,3}^{(1)} & q_{2,4}^{(1)} & q_{2,5}^{(1)} \end{bmatrix} \text{ and } \begin{bmatrix} q_{0,1}^{(1)} & q_{0,2}^{(1)} & q_{0,3}^{(1)} & q_{0,4}^{(1)} & q_{0,5}^{(1)} \end{bmatrix}$$

$$a^{(2)} = \begin{bmatrix} a_{0,0}^{(2)} & a_{0,1}^{(2)} & a_{0,2}^{(2)} \\ a_{1,0}^{(2)} & a_{1,1}^{(2)} & a_{1,2}^{(2)} \\ a_{2,0}^{(2)} & a_{2,1}^{(2)} & a_{2,2}^{(2)} \\ a_{3,0}^{(2)} & a_{3,1}^{(2)} & a_{3,2}^{(2)} \\ a_{4,0}^{(2)} & a_{4,1}^{(2)} & a_{4,2}^{(2)} \end{bmatrix} \rightarrow \begin{bmatrix} a_{1,0}^{(2)} & a_{1,1}^{(2)} & a_{1,2}^{(2)} \\ a_{2,0}^{(2)} & a_{2,1}^{(2)} & a_{2,2}^{(2)} \\ a_{3,0}^{(2)} & a_{3,1}^{(2)} & a_{3,2}^{(2)} \\ a_{4,0}^{(2)} & a_{4,1}^{(2)} & a_{4,2}^{(2)} \end{bmatrix} \text{ and } \begin{bmatrix} a_{0,0}^{(2)} & a_{0,1}^{(2)} & a_{0,2}^{(2)} \end{bmatrix}$$

So without column addition.

$$\text{original input } \begin{bmatrix} X_1^{(0)0} & X_2^{(0)0} & X_3^{(0)0} \\ X_1^{(0)1} & X_2^{(0)1} & X_3^{(0)1} \end{bmatrix} \cdot \begin{bmatrix} d_{1,1}^{(0)} & d_{1,2}^{(0)} \\ d_{2,1}^{(0)} & d_{2,2}^{(0)} \\ d_{3,1}^{(0)} & d_{3,2}^{(0)} \end{bmatrix}$$

$$\left\{ \begin{array}{l} X_1^{(0)0} \cdot q_{1,1}^{(0)} + X_2^{(0)0} \cdot q_{2,1}^{(0)} + X_3^{(0)0} \cdot q_{3,1}^{(0)} \\ X_1^{(0)1} \cdot q_{1,1}^{(0)} + X_2^{(0)1} \cdot q_{2,1}^{(0)} + X_3^{(0)1} \cdot q_{3,1}^{(0)} \end{array} \right. \quad \left. \begin{array}{l} X_1^{(0)0} \cdot q_{1,2}^{(0)} + X_2^{(0)0} \cdot q_{2,2}^{(0)} + X_3^{(0)0} \cdot q_{3,2}^{(0)} \\ X_1^{(0)1} \cdot q_{1,2}^{(0)} + X_2^{(0)1} \cdot q_{2,2}^{(0)} + X_3^{(0)1} \cdot q_{3,2}^{(0)} \end{array} \right\} \rightarrow \text{this operation lacks the bias addition.}$$

$$\left[ \begin{array}{ll} X_0^{(0)0} \cdot q_{0,1}^{(0)} & X_0^{(0)0} \cdot q_{0,2}^{(0)} \\ X_0^{(0)1} \cdot q_{0,1}^{(0)} & X_0^{(0)1} \cdot q_{0,2}^{(0)} \end{array} \right] \rightarrow \text{this matrix should be added to prev calculation.}$$

if can be represented

$$\begin{bmatrix} +1 \\ +1 \end{bmatrix} \cdot \begin{bmatrix} q_{0,1}^{(0)} & q_{0,2}^{(0)} \end{bmatrix}$$

$X^{(0)}$  are bias weights.

this matrix's size must be equal to input matrix's row number. If there are 5 samples it must be  $(5, 1)$

If  $10 \rightarrow (10, 1)$

$$\begin{bmatrix} +1 \\ +1 \\ +1 \\ +1 \\ +1 \end{bmatrix}$$

so for all weights it can be shown that separated matrix

$$O^{(1)} = X^{(0)} \cdot q^{(0)} + \begin{bmatrix} +1 \\ +1 \end{bmatrix} \cdot b^{(0)}$$

$$X^{(1)} = \bar{v}(O^{(1)})$$

$$O^{(2)} = X^{(1)} \cdot q^{(1)} + \begin{bmatrix} +1 \\ +1 \end{bmatrix} \cdot b^{(1)}$$

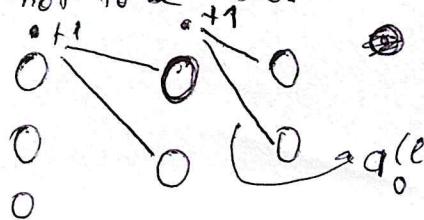
$$X^{(2)} = \bar{v}(O^{(2)})$$

$$O^{(3)} = X^{(2)} \cdot q^{(2)} + \begin{bmatrix} +1 \\ +1 \end{bmatrix} \cdot b^{(2)} \rightarrow \text{for regression}$$

it is output  
for classification softmax  $(O^{(3)})$ .

### Back Propagation

Previous backpropagation derivation is based on column vector addition.  
In this configuration, we have to ~~but~~  $a^{(2)}$  row vector because they represent the bias weights. while trasferring the error to previous layer these weights need not to be used.



There is no error backpropagation from these weights. So if we exclude these weights in calculation.

$$\Delta a^{(2)} = \eta \cdot X^{(2)T} \cdot Err$$

for ex.

$$Err = \begin{bmatrix} r_0^0 & r_1^0 & r_2^0 \\ r_0^1 & r_1^1 & r_2^1 \end{bmatrix}$$

$$X^{(2)} = \begin{bmatrix} +1 & X_1^{(2)0} & X_2^{(2)0} & X_3^{(2)0} & X_4^{(2)0} \\ +1 & X_1^{(2)1} & X_2^{(2)1} & X_3^{(2)1} & X_4^{(2)1} \end{bmatrix}$$

$$\Delta a^{(2)} = \eta \cdot \begin{bmatrix} +1 & +1 \\ X_1^{(2)0} & X_1^{(2)1} \\ X_2^{(2)0} & X_2^{(2)1} \\ X_3^{(2)0} & X_3^{(2)1} \\ X_4^{(2)0} & X_4^{(2)1} \end{bmatrix} \cdot \begin{bmatrix} r_0^0 & r_1^0 & r_2^0 \\ r_0^1 & r_1^1 & r_2^1 \end{bmatrix}$$

$$\begin{aligned} \Delta a_{0,0}^{(2)} &= \Delta a_{0,1}^{(2)} = \Delta a_{0,2}^{(2)} \\ \Delta a_{1,0}^{(2)} &= \Delta a_{1,1}^{(2)} = \Delta a_{1,2}^{(2)} \\ \Delta a_{2,0}^{(2)} &= \Delta a_{2,1}^{(2)} = \Delta a_{2,2}^{(2)} \\ \Delta a_{3,0}^{(2)} &= \Delta a_{3,1}^{(2)} = \Delta a_{3,2}^{(2)} \\ \Delta a_{4,0}^{(2)} &= \Delta a_{4,1}^{(2)} = \Delta a_{4,2}^{(2)} \end{aligned}$$

so to update  $b^{(0)}$   
we can get it from  
above

So similar to what we have done for forward backpropagation, while updating biases we need to use a row vector of ones.

$$\Delta a^{(1)} = \eta \cdot X^{(1)T} \cdot \text{Layer 2 error}$$

$$\eta [+1 \ 1] \cdot \begin{bmatrix} r_0^0 & r_1^0 & r_2^0 \\ r_0^1 & r_1^1 & r_2^1 \end{bmatrix}$$

$$\begin{bmatrix} \Delta a_{0,0}^{(2)} & \Delta a_{0,1}^{(2)} & \Delta a_{0,2}^{(2)} \end{bmatrix}$$

Previously

$$X^{(1)T} \rightarrow \eta \cdot \begin{bmatrix} +1 & +1 \\ X_1^{(1)0} & X_1^{(1)1} \\ X_2^{(1)0} & X_2^{(1)1} \end{bmatrix} \cdot \begin{bmatrix} r_0^0 & r_1^0 & r_2^0 \\ r_0^1 & r_1^1 & r_2^1 \\ r_1^0 & r_1^1 & r_2^0 \\ r_1^1 & r_2^0 & r_2^1 \\ r_2^0 & r_2^1 & r_2^0 \end{bmatrix}$$

$$\text{In new configuration } X^{(1)T} \rightarrow \begin{bmatrix} X_1^{(1)0} & X_1^{(1)1} \\ X_2^{(1)0} & X_2^{(1)1} \end{bmatrix} \cdot \begin{bmatrix} r_0^0 & r_1^0 & r_2^0 & r_3^0 \\ r_0^1 & r_1^1 & r_2^1 & r_3^1 \end{bmatrix} \rightarrow \begin{bmatrix} \Delta a_{0,1}^{(1)} & \Delta a_{0,2}^{(1)} & \Delta a_{0,3}^{(1)} & \Delta a_{0,4}^{(1)} \\ \Delta a_{1,1}^{(1)} & \Delta a_{1,2}^{(1)} & \Delta a_{1,3}^{(1)} & \Delta a_{1,4}^{(1)} \\ \Delta a_{2,1}^{(1)} & \Delta a_{2,2}^{(1)} & \Delta a_{2,3}^{(1)} & \Delta a_{2,4}^{(1)} \end{bmatrix}$$

$$\Delta a^{(1)} \text{ so for } \Delta b^{(1)} \rightarrow \eta [+1 \ 1] \cdot \text{layer 2 error}$$

so for changing bias weight, while updating bias these weights should be multiplied with a row of one. Its size should be as input matrix row count.

$$\Delta a^{(2)} = \eta \cdot X^{(2)T} \cdot \text{Output Error}$$

$$\Delta b^{(2)} = \eta [+1 \ 1 \dots] \cdot \text{Output Error}$$

$$\Delta a^{(1)} = \eta \cdot X^{(1)T} \cdot \text{Layer 2 Error}$$

$$\Delta b^{(1)} = \eta [+1 \ 1 \dots] \cdot \text{Layer 2 Error}$$

$$\Delta a^{(0)} = \eta \cdot X^{(0)T} \cdot \text{Layer 1 Error}$$

$$\Delta b^{(0)} = \eta [+1 \ 1 \dots] \cdot \text{Layer 1 Error}$$

$$\Delta a^{(e)} = \eta \cdot X^{(e)T} \cdot Err^{(e+1)}$$

$$\Delta b^{(e)} = \eta [+1 \ 1 \dots] \cdot Err^{(e+1)}$$