

$K \in \{0, 1, 2, \dots, m\}$ $K=0$: input layer
 $K=m$: output layer

$X \in \mathbb{R}^{N \times M}$ N : instances
 M : features

for $1 \leq K \leq m$:

$[W^K]_{ij} = W^K_{ij}$
 $i \in \{1, \dots, r_{k-1}\}$ where $r_0 = N$
 $j \in \{1, \dots, r_k\}$ $r_m = 1$

$$W^K = \begin{bmatrix} b^K_1 & b^K_2 & \dots & b^K_{r_k} \\ W^K_{11} & W^K_{12} & \dots & W^K_{1r_k} \\ \vdots & \vdots & \ddots & \vdots \\ W^K_{r_{k-1}1} & W^K_{r_{k-1}2} & \dots & W^K_{r_{k-1}r_k} \end{bmatrix}$$

$$\Theta^K = \begin{bmatrix} 1 & g(a^K_{11}) & g(a^K_{12}) & \dots & g(a^K_{1r_k}) \\ 1 & g(a^K_{21}) & g(a^K_{22}) & \dots & g(a^K_{2r_k}) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & g(a^K_{N1}) & g(a^K_{N2}) & \dots & g(a^K_{Nr_k}) \end{bmatrix}$$

$$A^K = \Theta^{K-1} W^K \Rightarrow A^K = \begin{bmatrix} a^K_{11} & a^K_{12} & \dots & a^K_{1r_k} \\ a^K_{21} & a^K_{22} & \dots & a^K_{2r_k} \\ \vdots & \vdots & \ddots & \vdots \\ a^K_{N1} & a^K_{N2} & \dots & a^K_{Nr_k} \end{bmatrix}$$

$$\Rightarrow \Theta^K = [\mathbf{1}, g(A^K)] \in \mathbb{R}^{N \times (r_k+1)}$$

$$\Theta^0 = \begin{bmatrix} 1 & X_{11} & X_{12} & \dots & X_{1M} \\ 1 & X_{21} & X_{22} & \dots & X_{2M} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{N1} & X_{N2} & \dots & X_{NM} \end{bmatrix}$$

$W^K \in \mathbb{R}^{(r_{k-1}+1) \times r_k}$

$\Theta^{K-1} \in \mathbb{R}^{N \times (r_{k-1}+1)}$

$A^K \in \mathbb{R}^{N \times r_k}$

Denote $W^K = \begin{bmatrix} \vec{b}^K \\ W^K_* \end{bmatrix}$ and

$\Theta^K = [\mathbf{1}, \Theta^K_*] \Rightarrow \Theta^K_* = g(A^K)$

Define $\Delta^m = \hat{Y} - Y \Rightarrow \Delta^m = \begin{bmatrix} \hat{y}_1 - y_1 \\ \hat{y}_2 - y_2 \\ \vdots \\ \hat{y}_N - y_N \end{bmatrix}$

$\hat{Y} = g(A^m)$

where $A^m \in \mathbb{R}^{N \times 1}$

$\Delta^K = \Theta^K_* \odot (\mathbf{I}_{(N \times r_k)} - \Theta^K_*) \odot [\Delta^{K+1} (W^{K+1}_*)^T]$

where $A \odot B = a_{ij} b_{ij}$
Hadamard product

make ~~transformation~~ dimensional transforms:

where $T_1(\Theta^{K-1}) \in \mathbb{R}^{N \times (r_{k-1}+1) \times 1}$ and $T_2(\Delta^K) \in \mathbb{R}^{N \times 1 \times r_k}$

Define $\partial \mathcal{E}^K = T_1(\Theta^{K-1}) \odot T_2(\Delta^K) \Rightarrow \partial \mathcal{E}^K \in \mathbb{R}^{N \times (r_{k-1}+1) \times r_k}$

$\Rightarrow \partial \mathcal{E}^K = [\partial \mathcal{E}^K_1, \partial \mathcal{E}^K_2, \dots, \partial \mathcal{E}^K_N]; DE^K = \frac{1}{N} \sum_{d=1}^N [\partial \mathcal{E}^K]_d$

where $[\partial \mathcal{E}^K]_d \in \mathbb{R}^{(r_{k-1}+1) \times r_k}$ and $DE^K \in \mathbb{R}^{(r_{k-1}+1) \times r_k}$

iterate till $\hat{W}^K = W^K - \alpha DE^K$ converges

0 0 0
0 0 0 0
0 0 0

Given: $[X, Y]$
 $X \in \mathbb{R}^{N \times M}$

Set: W^k for $1 \leq k \leq m$

$r_0 = N, r_m = 1$
decide r_k for $1 \leq k \leq m-1$

W^k for $1 \leq k \leq m, W^k \in \mathbb{R}^{(r_{k-1}) \times r_k}$

k in range $(1, m)$

algorithm

$k \in \{0, 1, 2, 3\}$

$$\Theta^0 = [I, X]$$

for $1 \leq k \leq m-1$:

$$A^k = \Theta^{k-1} W^k$$

$$\Theta_*^k = \sigma(A^k)$$

$$\Theta^k = [I, \Theta_*^k]$$

$$A^m = \Theta^{m-1} W^m$$

$$\Delta^m = \hat{Y} - Y$$

for $m-1 \geq k \geq 1$:

$$\Delta^k = \Theta_*^k \odot (I_{(W^k r_k)} - \Theta_*^k) \odot [\Delta^{k+1} (W_*^{k+1})^T]$$

for $1 \leq k \leq m-1$:

$$\Theta^{k-1} \mapsto T_1(\Theta^{k-1})$$

$$\Delta^k \mapsto T_2(\Delta^k)$$

$$\partial \mathcal{E}^k = T_1(\Theta^{k-1}) \odot T_2(\Delta^k)$$

$$D\mathcal{E}^k = \frac{1}{N} \sum_{d=1}^N [\partial \mathcal{E}^k]_d$$

$$\hat{W}^k = W^k - \alpha D\mathcal{E}^k$$