

$$\begin{aligned} \underline{X} \cdot \underline{W}^{(1)} + \underline{b}_1 &= \underline{Z} \\ \underline{Z} \cdot \underline{W}^{(2)} + \underline{b}_2 &= \underline{\hat{Y}} \end{aligned} \quad \left| \quad \begin{aligned} \underline{b}_1 &= \begin{bmatrix} b_1 \\ \vdots \\ b_4 \end{bmatrix} \}^n \\ \underline{b}_2 &= \begin{bmatrix} b_2 \\ \vdots \\ b_1 \end{bmatrix} \}^n \end{aligned} \right.$$

scalar loss = MSE = $\frac{\sum_i (\hat{y}_i - y_i)^2}{n} = \|\hat{\underline{Y}} - \underline{Y}\|_2^2 \cdot \frac{1}{n} = (\hat{\underline{Y}} - \underline{Y})^T (\hat{\underline{Y}} - \underline{Y}) \cdot \frac{1}{n}$

①

$$\begin{aligned} 1 \times 1 \quad \text{gradient } \hat{y}_n &= \nabla_{\hat{y}_n} \text{loss} = \frac{2}{n} (\hat{y}_n - y_n) \\ 1 \times 1 \quad \text{gradient } y_n &= \nabla_{y_n} \text{loss} = -\frac{2}{n} (\hat{y}_n - y_n) \end{aligned} \quad \left. \vphantom{\begin{aligned} 1 \times 1 \quad \text{gradient } \hat{y}_n \\ 1 \times 1 \quad \text{gradient } y_n \end{aligned}} \right\} \text{single row}$$

$$\begin{aligned} \Rightarrow \text{gradient } \underline{\hat{Y}} &= \nabla_{\underline{\hat{Y}}} \text{loss} = \nabla_{\underline{\hat{Y}}} \frac{\underline{\hat{Y}}^T \underline{\hat{Y}} - \underline{\hat{Y}}^T \underline{Y} - \underline{Y}^T \underline{\hat{Y}} + \underline{Y}^T \underline{Y}}{n} = \left(\frac{\underline{\hat{Y}}}{n} + \frac{\underline{\hat{Y}}}{n} \right) - \frac{\underline{Y}}{n} - \frac{\underline{Y}}{n} + 0 \\ \text{gradient } \underline{Y} &= \nabla_{\underline{Y}} \text{loss} = \frac{-\underline{\hat{Y}} - \underline{\hat{Y}} + \underline{Y} + \underline{Y}}{n} = \frac{2}{n} (\underline{Y} - \underline{\hat{Y}}) = -\frac{2}{n} (\underline{\hat{Y}} - \underline{Y}) \end{aligned}$$

for each row in $\hat{\underline{y}}$:

$$\begin{cases} s_{n1} \cdot w_1^{(2)} + s_{n2} \cdot w_2^{(2)} + s_{n3} \cdot w_3^{(2)} + s_{n4} \cdot w_4^{(2)} + b_2 = \hat{y}_n \\ \vdots \\ s_{n1} \cdot w_1^{(2)} + s_{n2} \cdot w_2^{(2)} + s_{n3} \cdot w_3^{(2)} + s_{n4} \cdot w_4^{(2)} + b_2 = \hat{y}_n \end{cases}$$

$$= \begin{cases} \underline{s}_1 \cdot \underline{w}^{(2)} + b_2 = \hat{y}_1 \\ \vdots \\ \underline{s}_n \cdot \underline{w}^{(2)} + b_2 = \hat{y}_n \end{cases}$$

②

$$\begin{aligned} 1 \times 4 \quad \text{gradient } \underline{s}_n &= \text{gradient } \hat{y}_n \cdot \nabla_{\underline{s}_n} \hat{y}_n = \frac{2}{n} (\hat{y}_n - y_n) \cdot \underline{w}^{(2)T} \\ 4 \times 1 \quad \text{gradient } \underline{w}^{(2)} &= \nabla_{\underline{w}^{(2)}} \hat{y}_n \cdot \text{gradient } \hat{y}_n = \underline{s}_1^T \cdot \frac{2}{n} (\hat{y}_n - y_n) \\ 1 \times 1 \quad \text{gradient } b_2 &= \nabla_{b_2} \hat{y}_n \cdot \text{gradient } \hat{y}_n = 1 \cdot \frac{2}{n} (\hat{y}_n - y_n) \end{aligned} \quad \left. \begin{array}{l} \text{single} \\ \text{row} \end{array} \right\}$$

$$\Rightarrow \begin{aligned} n \times 4 \quad \text{gradient } \underline{s} &= \begin{bmatrix} \text{gradient } s_1 \\ \vdots \\ \text{gradient } s_n \end{bmatrix} = \frac{2}{n} \begin{bmatrix} \hat{y}_1 - y_1 \\ \vdots \\ \hat{y}_n - y_n \end{bmatrix} \cdot \underline{w}^{(2)T} = -\frac{2}{n} (y - \hat{y}) \cdot \underline{w}^{(2)T} = \text{gradient } \hat{\underline{y}} \cdot \underline{w}^{(2)T} \\ 4 \times 1 \quad \text{gradient } \underline{w}^{(2)} &= \nabla_{\underline{w}^{(2)}} \hat{\underline{y}} \cdot \text{gradient } \hat{\underline{y}} = \begin{bmatrix} \underline{s}_1^T \\ \vdots \\ \underline{s}_n^T \end{bmatrix} \cdot -\frac{2}{n} (y - \hat{y}) = \underline{s}^T \cdot -\frac{2}{n} (y - \hat{y}) \\ 1 \times 1 \quad \text{gradient } b_2 &= \nabla_{b_2} \hat{\underline{y}} \cdot \text{gradient } \hat{\underline{y}} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \cdot -\frac{2}{n} (y - \hat{y}) = -\frac{2}{n} \cdot \sum_{i=1}^n (y_i - \hat{y}_i) \end{aligned}$$

scalar

③

$$\begin{matrix} n \times 3 & 3 \times 4 & n \times 4 & & n \times 4 \\ \underline{X} & \cdot \underline{W}^{(1)} & + \underline{b}_1 & = & \underline{Z} \end{matrix} \quad \left\{ \begin{matrix} 1 \times 3 & 3 \times 4 & & 1 \times 4 & 1 \times 4 \\ \underline{x}_1 & \underline{W}^{(1)} & + & \underline{b}_1 & = \underline{z}_1 \\ & \vdots & & \vdots & \\ & \underline{x}_n & \underline{W}^{(1)} & + & \underline{b}_1 = \underline{z}_n \end{matrix} \right.$$

$$\text{gradient } x_{ij} = \sum_{k=1}^4 \text{gradient } S_{ik} \cdot \nabla_{x_{ij}} S_{ik} = \sum_{k=1}^4 -\frac{2}{n} (y_i - \hat{y}_i) \cdot w_k^{(2)} \cdot w_{jk}^{(1)} \\ = -\frac{2}{n} (y_i - \hat{y}_i) \cdot \underline{w}^{(2)T} \cdot \underline{w}_j^{(1)T}$$

$$\Rightarrow \text{gradient } \underline{x}_n = -\frac{2}{n} (y_i - \hat{y}_i) \cdot \underline{w}^{(2)T} \cdot \underline{w}^{(1)T} = \text{gradient } \underline{z}_n \cdot \underline{w}^{(1)T}$$

$$\Rightarrow \text{gradient } \underline{X} = -\frac{2}{n} (Y - \hat{Y}) \cdot \underline{w}^{(2)T} \cdot \underline{w}^{(1)T} = \text{gradient } \underline{Z} \cdot \underline{w}^{(1)T}$$

$$\text{gradient } w_{ij}^{(1)} = \sum_{k=1}^n \text{gradient } S_{kj} \cdot \nabla_{w_{ij}^{(1)}} S_{kj} = \sum_{k=1}^n -\frac{2}{n} (y_k - \hat{y}_k) \cdot w_j^{(2)} \cdot x_{ki} \\ = -\frac{2}{n} w_j^{(2)} \cdot \begin{bmatrix} x_{1i} \\ \vdots \\ x_{ni} \end{bmatrix}^T \cdot (Y - \hat{Y})$$

$$\Rightarrow \text{gradient } \underline{w}_n^{(1)} = -\frac{2}{n} \underline{w}^{(2)T} \cdot \begin{bmatrix} x_{1i} \\ \vdots \\ x_{ni} \end{bmatrix}^T \cdot (Y - \hat{Y})$$

1x4 1x4 scalar

$$\Rightarrow \text{gradient } \underline{W}^{(1)} = -\frac{2}{n} \begin{bmatrix} x_{11} & \dots & x_{n1} \\ \vdots & & \vdots \\ x_{13} & \dots & x_{n3} \end{bmatrix} \cdot (Y - \hat{Y}) \cdot \underline{w}^{(2)T} = -\frac{2}{n} \underline{X}^T \cdot (Y - \hat{Y}) \cdot \underline{w}^{(2)T} = \underline{X}^T \cdot \text{gradient } \underline{Z}$$

3x4 3xn nx1 1x4

$$\text{gradient } b_{1i} = \sum_{k=1}^n \text{gradient } S_{ki} \cdot \nabla_{b_{1i}} S_{ki} = \sum_{k=1}^n -\frac{2}{n} (y_k - \hat{y}_k) \cdot w_i^{(2)} \cdot 1$$

$$= -\frac{2}{n} \cdot \sum_{k=1}^n (y_k - \hat{y}_k) \cdot w_i^{(2)}$$

$$\Rightarrow \text{gradient } \underline{b}_1 = -\frac{2}{n} \cdot \sum_{k=1}^n (y_k - \hat{y}_k) \cdot \underline{w}^{(2)T}$$

$1 \times 4 \qquad \qquad \qquad 1 \times 4$

Brief Summary

for linear node: $\underline{X} \cdot \underline{W} + \underline{b} = \underline{S}$

$n \times d \quad d \times c \quad n \times c \qquad \qquad n \times c$



$$\text{gradient } \underline{X} = \text{gradient } \underline{S} \cdot \underline{W}^T$$

$n \times d \qquad \qquad n \times c \qquad \qquad c \times d$

$$\text{gradient } \underline{W} = \underline{X}^T \cdot \text{gradient } \underline{S}$$

$d \times c \qquad \qquad d \times n \qquad \qquad n \times c$

$$\text{gradient } \underline{b} = \text{sum}(\text{gradient } \underline{S}, \text{axis}=0)$$

for activation Node: $\Delta(\underline{X}) = \underline{Y}$

$m \times n \qquad \qquad m \times n$

$$\text{Sigmoid } \Delta(x) = \frac{1}{1 + e^{-x}}$$

$$\text{Tan h } \Delta(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$\text{gradient } \Delta(x) = \text{gradient } \underline{Y} \cdot \Delta(1 - \Delta)$$

$$\text{gradient } \Delta(x) = \text{gradient } \underline{Y} \cdot (1 - \tanh^2(x))$$

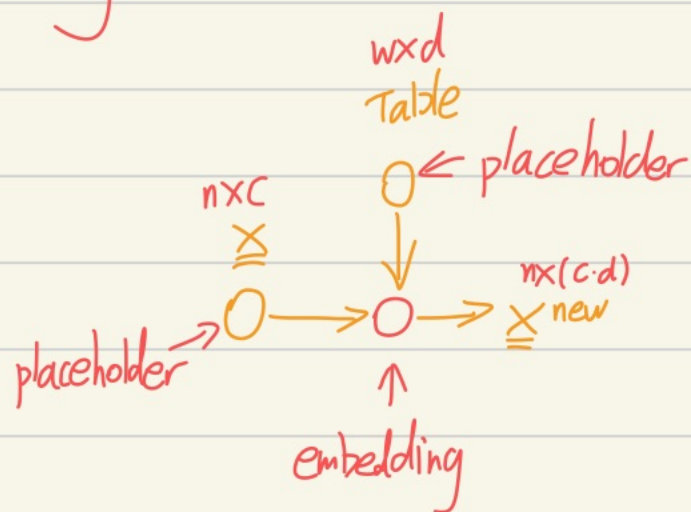
for Add Node: $\underline{X} + \underline{Y} + \underline{b} = \underline{Z}$

$\underline{b} = \begin{bmatrix} b \\ \vdots \\ b \end{bmatrix}$ (size $n \times 1$)

$\text{gradient } \underline{X} = \text{gradient } \underline{Y} = \text{gradient } \underline{Z}$

$\text{gradient } \underline{b} = \text{Sum}(\text{gradient } \underline{Z}, \text{axis}=0)$

for Embedding case in NLP



\underline{X} : n sentences, c words in each

Table: Total w distinct words in all sentences, word embedding d dimension

In embedding node:

① change $\underline{X} \rightarrow \underline{\pi}$ [one hot transformation]

② $\underline{\pi} \cdot \text{Table} = \underline{R} \xrightarrow{\text{reshape}} \underline{X}_{\text{new}}$

$$\text{gradient } \underline{R} = \text{reshape } \text{gradient } \underline{x}^{\text{new}} \text{ to } nc \times d$$

$(nc) \times d$ $n \times (c \cdot d)$

$$\text{gradient Table} = \underline{x}^T \cdot \text{gradient } \underline{R}$$

$w \times d$ $w \times (nc)$ $(nc) \times d$

$$\text{gradient } \underline{z} = \text{gradient } \underline{R} \cdot \text{Table}^T$$

$(nc) \times w$ $(nc) \times d$ $d \times w$

$$\text{gradient } \underline{x} = \text{reshape sum}(\text{gradient } \underline{z}, \text{axis} = 1) \text{ to } n \times c$$

$n \times c$ $h \times c$ $nc \times w$