**1.**

**1st:**



$b_1$ $^{1\times H}$

$X$ $^{N\times D}$

$W_1$ $^{D\times H}$

$(\ast) \to (+) \xrightarrow{z} (ReLu) \to H$ $^{N\times H}$

$R(W_1)$

**2nd:**

$b_2$ $^{1\times C}$

$H$ $^{N\times H}$

$W_2$ $^{H\times C}$

$(\ast) \to (+) \xrightarrow{S \;^{N\times C}} (softmax) \to (+) \to L$

$R(W_2)$

---

**2.**

**2nd:** $\frac{\partial L}{\partial H}$ $H$

$b_2$ $\frac{\partial L}{\partial b_2}$

$\frac{\partial L}{\partial W_2}$ $W_2$

$(\ast) \to (+) \xrightarrow{S,\; \frac{\partial L}{\partial S}} (softmax) \to (+) \to L$

$R(W_2)$ $\frac{\partial R}{\partial W_2}$

$$L = \frac{1}{N}\sum_{i=1}^{N} L_i + \lambda \cdot R(W) \;;$$
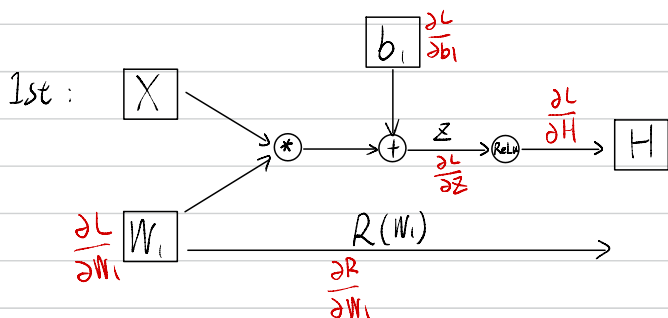
$$R(W) = \left(\sum W_1^2 + \sum W_2^2\right)$$

$$\frac{\partial R}{\partial W_2} = 2 \times W_2$$

$$\frac{\partial L}{\partial S} = \begin{cases} \frac{1}{N}(P_{i,j}-1), & j=y_i \\ \frac{1}{N}(P_{i,j}), & j\neq y_i \end{cases} \in R^{N\times C}$$

$$\frac{\partial L}{\partial b_2} = 1_N^T \frac{\partial L}{\partial S} \in R^{1\times C}$$

$$\frac{\partial L}{\partial H} = \frac{\partial L}{\partial S} W_2^T \in R^{N\times H}$$

$$\frac{\partial L}{\partial W_2} = H^T \frac{\partial L}{\partial S} + 2\cdot\lambda W_2 \in R^{H\times C}$$

---

**1st:**

$b_1$ $\frac{\partial L}{\partial b_1}$

$X$

$\frac{\partial L}{\partial W_1}$ $W_1$

$(\ast) \to (+) \xrightarrow{z,\; \frac{\partial L}{\partial z}} (ReLu) \xrightarrow{\frac{\partial L}{\partial H}} H$

$R(W_1)$ $\frac{\partial R}{\partial W_1}$

$$\frac{\partial L}{\partial z} = \begin{cases} \left[\frac{\partial L}{\partial H}\right]_{i,j}, & z_{i,j}>0 \\ 0, & \text{otherwise} \end{cases} \in R^{N\times H}$$

$$\frac{\partial L}{\partial b_1} = 1_N^T \cdot \frac{\partial L}{\partial z} \in R^{1\times H}$$

$$\frac{\partial R}{\partial W_1} = 2 W_1$$

$$\frac{\partial L}{\partial W_1} = X^T \frac{\partial L}{\partial z} + 2\lambda W_1 \in R^{D\times H}$$

**3.**

$$g_t = \nabla L_t$$

$$M_t = \beta \cdot M_{t-1} + (1-\beta) g_t$$

$$M_0 = 0$$

$$M_1 = \beta \cdot M_0 + (1-\beta) \cdot g_1$$

$$= (1-\beta) \cdot g_1$$

\* $\hat{m}_1 = \dfrac{M_1}{1-\beta}$

$$M_2 = \beta \cdot M_1 + (1-\beta) g_2$$

$$= \beta(1-\beta) g_1 + (1-\beta) g_2$$

$$M_3 = \beta \cdot M_2 + (1-\beta) g_3$$

$$= \beta^2 (1-\beta) g_1 + \beta(1-\beta) g_2 + (1-\beta) g_3$$

$$\vdots$$

$$M_t = (1-\beta) \sum_{i=1}^{t} \beta^{t-i} \cdot g_i$$

$$E[M_t] = (1-\beta) \sum_{i=1}^{t} \beta^{t-i} \cdot E[g_i]$$

Assume $E[g_i] = E[g]$

$$E[M_t] = (1-\beta) \sum_{i=1}^{t} \beta^{t-i} \cdot E[g]$$

$$= (1-\beta) \cdot \frac{1-\beta^t}{1-\beta} E[g]$$

$$= (1-\beta^t) \cdot E[g]$$

$$\hat{m}_t = \frac{M_t}{1-\beta^t}$$