**Model :**



$X^{(0)}$  $W_1$  $X^{(1)}$  $W_2$  $X^{(2)}$

**Detailed :**



$X^{(0)}$  $W^{(1)}$  $z^{(1)}$  $\varphi$  $X^{(1)}$  $W^{(2)}$  $z^{(2)}$  $\varphi$  $X^{(2)}$

$b^{(1)}$  $b^{(2)}$

**Params :**

$$W^{(1)} \in \mathbb{R}^{K \times 784}$$

$$W^{(2)} \in \mathbb{R}^{K \times 10}$$

$$b^{(1)} \in \mathbb{R}^{K}$$

$$b^{(2)} \in \mathbb{R}^{10}$$

**Loss :**

$$\mathcal{L}_n(x^{(n)}) = \frac{1}{2}\left( \sum_{k=1}^{10} \left[ f_k(x^{(n)}) - y_k^{(n)} \right]^2 \right)$$

**Goal :**

$$\frac{\partial \mathcal{L}_n}{\partial w_{ij}^{(1)}} \qquad \frac{\partial \mathcal{L}_n}{\partial b_i^{(1)}}$$
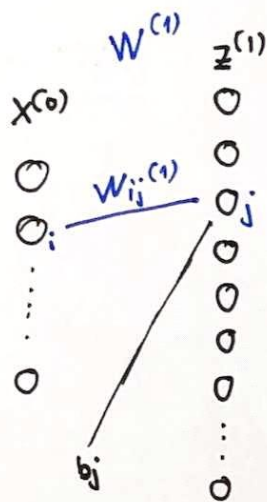
$$\frac{\partial \mathcal{L}_n}{\partial w_{ij}^{(2)}} \qquad \frac{\partial \mathcal{L}_n}{\partial b_i^{(2)}}$$

**Notation :** Let $\partial_j^{(1)} = \dfrac{\partial \mathcal{L}_n}{\partial z_j^{(1)}}$

$$\partial_j^{(2)} = \frac{\partial \mathcal{L}_n}{\partial z_j^{(2)}}$$

$$\partial_j^{(3)} = \frac{\partial \mathcal{L}_n}{\partial \ ?}$$

First Layer :



By the Chain Rule;

$$\frac{\partial \mathcal{L}_n}{\partial w_{ij}^{(1)}} = \sum_{k=1}^{K} \frac{\partial L^n}{\partial z_k^{(1)}} \cdot \frac{\partial z_k^{(1)}}{\partial w_{ij}^{(1)}}$$

$$= \frac{\partial \mathcal{Y}_n}{\partial z_j^{(1)}} \cdot \frac{\partial z_j^{(1)}}{\partial w_{ij}^{(1)}}$$

$$= \partial_j^{(1)} \cdot \frac{\partial z_j^{(1)}}{\partial w_{ij}^{(1)}}$$

Since $z^{(1)} = \left[W^{(1)}\right]^T x^{(0)} + b^{(1)}$,

$$z_j^{(1)} = \sum_{k=1}^{D} \left[W^{(1)}\right]_{jk}^T x_k^{(0)} + b_j = \sum_{k=1}^{D} w_{kj}^{(1)} x_k^{(0)} + b_j$$

Now $\dfrac{\partial z_j^{(1)}}{\partial w_{ij}^{(1)}} = \dfrac{\partial}{\partial w_{ij}^{(1)}} \left( \sum_{k=1}^{D} w_{kj}^{(1)} x_k^{(0)} + b_j \right)$

$$= \sum_{k=1}^{D} \frac{\partial}{\partial w_{ij}^{(1)}} \left( w_{kj}^{(1)} x_k^{(0)} \right)$$

$$= x_i^{(0)}$$

Hence $\dfrac{\partial \mathcal{L}_n}{\partial w_{ij}^{(1)}} = \partial_j^{(1)} \cdot x_i^{(0)}$ .

By the Chain Rule,
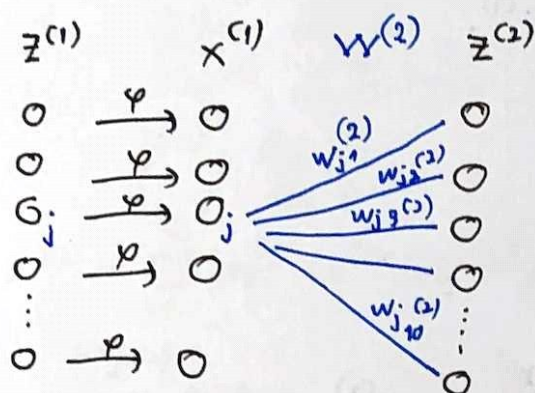
$$\frac{\partial \mathcal{Y}_n}{\partial b_i^{(1)}} = \sum_{k=1}^{D} \frac{\partial \mathcal{L}_n}{\partial z_k^{(1)}} \cdot \frac{\partial z_k^{(1)}}{\partial b_i^{(1)}}$$

$$= \frac{\partial \mathcal{L}_n}{\partial z_i^{(1)}} \cdot \frac{\partial z_i^{(1)}}{\partial b_i^{(1)}}$$

Since $z_i^{(1)} = \sum_{k=1}^{D} w_{ki}^{(1)} x_k^{(0)} + b_i$, $\dfrac{\partial z_i^{(1)}}{\partial b_i^{(1)}} = 1$.

Hence $\dfrac{\partial \mathcal{L}_n}{\partial b_i^{(1)}} = \dfrac{\partial \mathcal{Y}_n}{\partial z_i^{(1)}} = \partial_i^{(1)}$ .

Now we focus on $\theta_j^{(1)} = \dfrac{\partial \mathcal{L}_n}{\partial z_j^{(1)}}$.

$$z^{(1)} \qquad x^{(1)} \qquad W^{(2)} \qquad z^{(2)}$$



We have, by the Chain Rule,

$$\theta_j^{(1)} = \frac{\partial \mathcal{L}_n}{\partial z_j^{(1)}} = \frac{\partial \mathcal{L}_n}{\partial x_j^{(1)}} \cdot \frac{\partial x_j^{(1)}}{\partial z_j^{(1)}}$$

$$= \left( \sum_{k=1}^{10} \frac{\partial \mathcal{L}_n}{\partial z_k^{(2)}} \cdot \frac{\partial z_k^{(2)}}{\partial x_j^{(1)}} \right) \cdot \frac{\partial x_j^{(1)}}{\partial z_j^{(1)}}$$

$$= \left[ \sum_{k=1}^{10} \partial_k^{(2)} \cdot \frac{\partial z_k^{(2)}}{\partial x_j^{(1)}} \right] \cdot \frac{\partial x_j^{(1)}}{\partial z_j^{(1)}}$$

We have $x_j^{(1)} = \varphi(z_j^{(1)}) \implies \dfrac{\partial x_j^{(1)}}{\partial z_j^{(1)}} = \varphi'(z_j^{(1)})$

We have $z_k^{(2)} = \left( W^{(2)T} x^{(1)} + b^{(2)} \right)_k$

$$= \sum_{i=1}^{K} (w^2)_{ki}^T x_i^{(1)} + b_k^{(2)} = \sum_{i=1}^{K} w_{ik}^{(2)} x_i^{(1)} + b_k^{(2)}$$

Hence $\dfrac{\partial z_k^{(2)}}{\partial x_j^{(1)}} = \sum_{i=1}^{K} \dfrac{\partial}{\partial x_j^{(1)}} \left[ w_{ik}^{(2)} x_i^{(1)} \right] + b_k^{(2)}$

$$= w_{jk}^{(2)}$$

Now $\partial_j^{(1)} = \left( \sum_{k=1}^{10} \partial_k^{(2)} \cdot w_{jk} \right) \varphi'(z_j^{(1)})$

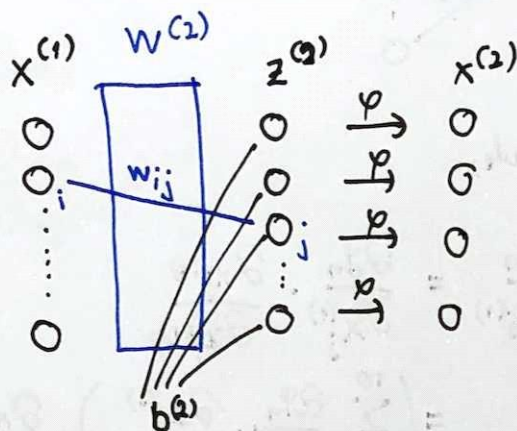$$= \left( \sum_{k=1}^{10} w_{jk}^{(2)} \partial_k^{(2)} \right) \varphi'(z_j^{(1)})$$

Observe that $\partial^{(1)} = W^{(2)} \partial^{(2)} \odot \varphi'(z^{(1)})$.

**Completing First Layer:**

$$\begin{cases} \dfrac{\partial \mathcal{L}_n}{\partial w_{ij}^{(1)}} = \partial_j^{(1)} \cdot x_i^{(0)} \\[2mm] \dfrac{\partial \mathcal{L}_n}{\partial b_i^{(1)}} = \partial_i^{(1)}. \end{cases}$$

Now consider the Second Layer:

## Second Layer



We want $\dfrac{\partial \mathcal{L}_n}{\partial w_{ij}^{(2)}}, \dfrac{\partial \mathcal{L}_n}{\partial b_i^{(2)}} = ?$

Now $\dfrac{\partial \mathcal{L}_n}{\partial w_{ij}^{(2)}} = \displaystyle\sum_{k=1}^{10} \dfrac{\partial z_k^{(2)}}{\partial w_{ij}^{(2)}} \cdot \dfrac{\partial \mathcal{L}^n}{\partial z_k^{(2)}}$

$= \displaystyle\sum_{k=1}^{10} \partial_k^{(2)} \cdot \dfrac{\partial z_k^{(2)}}{\partial w_{ij}^{(2)}}$

$= \partial_j^{(2)} \cdot \dfrac{\partial z_j^{(2)}}{\partial w_{ij}^{(2)}}$

We have $z_j^{(2)} = \left( \left( W^{(2)^T} x^{(1)} \right) + b^{(2)} \right)_j$

$= \displaystyle\sum_{k=1}^{K} W_{jk}^{(2)^T} \cdot x_k^{(1)} + b_j^{(2)}$

$= \displaystyle\sum_{k=1}^{K} W_{kj}^{(2)} x_k^{(1)} + b_j^{(2)}.$

Now $\dfrac{\partial z_j^{(2)}}{\partial w_{ij}^{(2)}} = \displaystyle\sum_{k=1}^{K} \dfrac{\partial}{\partial w_{ij}^{(2)}} \left( W_{kj}^{(2)} x_k^{(1)} \right)$

$= x_i^{(1)}$

Hence $\dfrac{\partial \mathcal{L}_n}{\partial w_{ij}^{(2)}} = \partial_j^{(2)} x_i^{(1)}$

We have
$$\frac{\partial \mathcal{L}_n}{\partial b_i^{(2)}} = \sum_{k=1}^{k} \frac{\partial z_k^{(2)}}{\partial b_i^{(2)}} \cdot \frac{\partial \mathcal{L}_n}{\partial z_k^{(2)}}$$

$$= \frac{\partial \mathcal{L}_n}{\partial z_i^{(2)}} \cdot \frac{\partial z_i^{(2)}}{\partial b_i^{(2)}}$$

$$= \partial_i^{(2)} \cdot \frac{\partial z_i^{(2)}}{\partial b_i^{(2)}}.$$

Since $z_i^{(2)} = \sum_{k=1}^{k} w_{ki}^{(2)} x_k^{(1)} + b_i^{(2)}$,

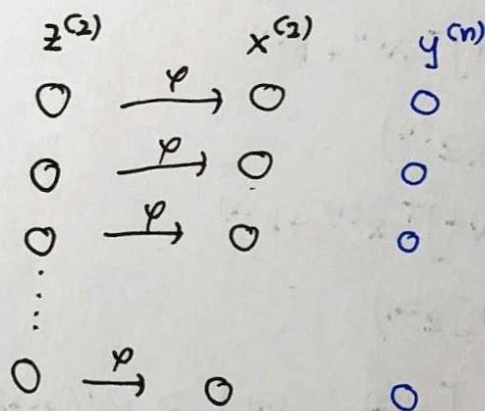$$\frac{\partial z_i^{(2)}}{\partial b_i^{(2)}} = 1.$$

Therefore, $\dfrac{\partial \mathcal{L}_n}{\partial b_i^{(2)}} = \partial_i^{(2)} \cdot 1 = \partial_i^{(2)}.$

It remains to compute $\partial_j^{(2)}$.

By definition, $\partial_j^{(2)} = \dfrac{\partial \mathcal{L}_n}{\partial z_j^{(2)}}.$

We consider the relevant segment of the network:



By the Chain Rule,

$$\partial_j^{(2)} = \frac{\partial \mathcal{L}_n}{\partial z_j^{(2)}} = \frac{\partial \mathcal{L}_n}{\partial x_j^{(2)}} \cdot \frac{\partial x_j^{(2)}}{\partial z_j^{(2)}}$$

We have $x_j^{(2)} = \varphi\left( z_j^{(2)} \right).$

Here $\dfrac{\partial x_j^{(2)}}{\partial z_j^{(2)}} = \varphi'\left( z_j^{(2)} \right).$

The interesting part is $\frac{\partial y^n}{\partial x_j^{(2)}}$.

We have $\quad \mathcal{L}_n(x^{(n)}) = \frac{1}{2} \sum_{k=1}^{10} (x_k^{(2)} - y_k^{(n)})^2$

Hence $\quad \frac{\partial \mathcal{L}_n}{\partial x_k^{(2)}} = \frac{1}{2} \cdot 2 \, (x_k^{(2)} - y_k^{(n)})$

$$= x_k^{(2)} - y_k^{(n)}$$

$$= \left[ x^{(2)} - y^{(n)} \right]_k.$$

Now $\quad \partial_j^{(2)} = \left[ x^{(2)} - y^{(n)} \right]_j \, \varphi'(z_j^{(2)})$

In vector form, $\quad \partial^{(2)} = (x^{(2)} - y^{(n)}) \odot \varphi'(z^{(2)})$.

Hence we are ready to state forward and backward pass.

We have forward:

$$x^{(0)} = x_n$$

$$z^{(1)} = [w^1]^T x^{(0)} + b^{(1)}$$

$$x^{(1)} = \phi(z^{(1)})$$

$$z^{(2)} = [w^2]^T x^{(1)} + b^{(2)}$$

$$x^{(2)} = \phi(z^{(2)})$$

$K \times 1 \times 1 \times 10$

Backward: $\quad \partial^{(2)} = \left[ x^{(2)} - y^{(n)} \right] \odot \varphi'(z^{(2)})$

$$\partial^{(1)} = w^{(2)} \partial^{(2)} \odot \varphi'(z^{(1)})$$

$$\begin{bmatrix} x_1^{(1)} \\ x_k^{(1)} \end{bmatrix} \left[ \partial_1^{(2)} \, \partial_2^{(2)} \dots \partial_{10}^{(2)} \right)$$

$\frac{\partial \mathcal{L}_n}{\partial w_{ij}^{(2)}} = \partial_j^{(2)} x_i^{(1)} \implies \frac{\partial \mathcal{L}_n}{\partial w^{(2)}} = x^{(1)} (\partial^2)^T \quad x_1 \partial_1 \; x_1 \partial_2 \; x_1 \partial_3$

$\frac{\partial \mathcal{L}_n}{\partial b_j^{(2)}} = \partial_j^{(2)} \quad \longrightarrow \quad \underline{\frac{\partial \mathcal{L}_n}{\partial b^{(2)}} = \partial^{(2)}}$

$$\underline{\frac{\partial \mathcal{L}_n}{\partial w^{(1)}} = x^{(0)} (\partial^1)^T}$$

$$\underline{\frac{\partial \mathcal{L}_n}{\partial b^{(1)}} = \partial^{(1)}}$$