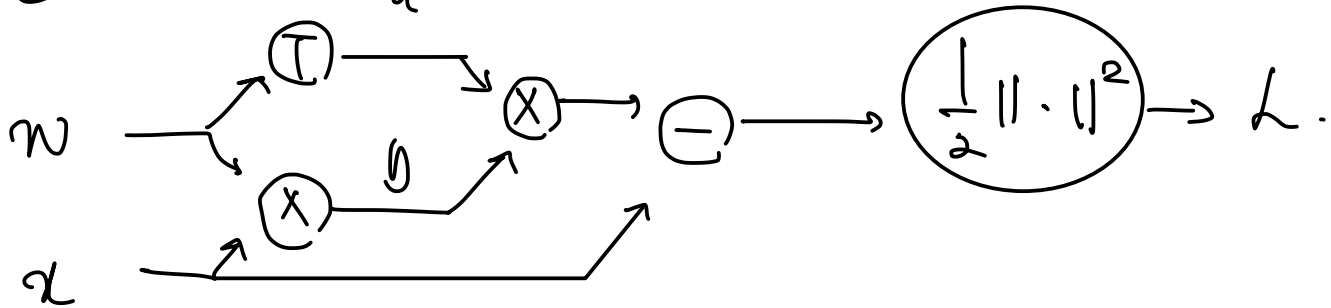# #1.

$$L = \frac{1}{2} \| W^T W x - x \|^2$$

(a)

The transformation $W^T W x$ first encodes information through $Wx$ and then decodes it through $W^T$.

The loss function $L$ measures how well the recovered information $W^T W x$ preserves the original input $x$. By minimizing the loss, the model learns to retain important features of $x$ during encoding and decoding.
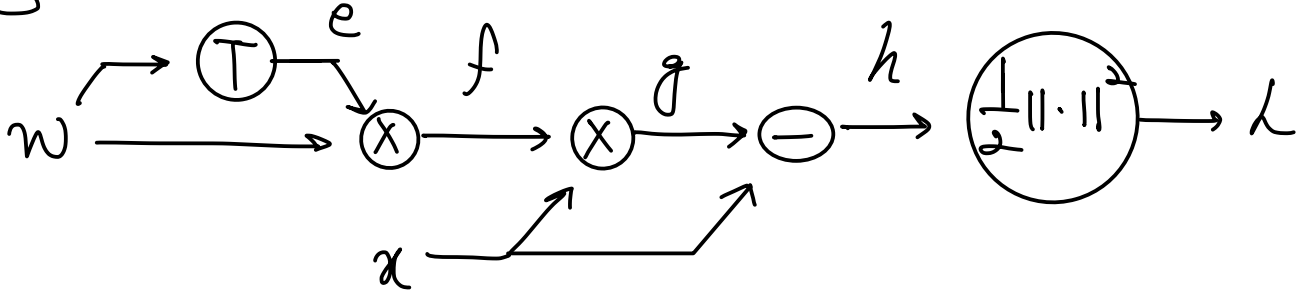
**(b)**

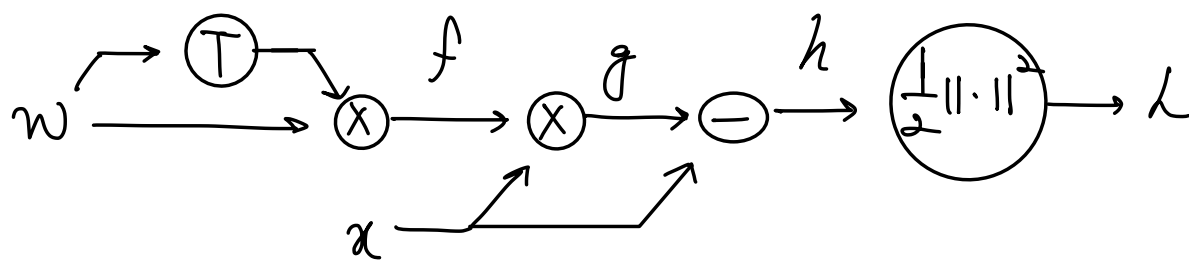$$\mathcal{L} = \frac{1}{2}\|w^T w x - x\|^2$$

① 



(OR)

② 



**(c)**

On the graph (b)-1, W contributes to two paths, $a$ and $b$
So, its gradient contributions must be accumulated accordingly
by the total derivative rule such that

$$\nabla_w \mathcal{L} = \frac{d\mathcal{h}}{da} \cdot \frac{da}{dw} + \frac{d\mathcal{L}}{db} \cdot \frac{db}{dw}$$

(d) $\quad L = \frac{1}{2}\|W^T W x - x\|^2$



$x \in \mathbb{R}^n$ $\qquad e \in \mathbb{R}^{n \times m}$ $\quad g \in \mathbb{R}^n$ $\quad L \in \mathbb{R}$

$W \in \mathbb{R}^{m \times n}$ $\qquad f \in \mathbb{R}^{n \times n}$ $\quad h \in \mathbb{R}^n$

---

$f = W^T W$ $\qquad\qquad L = \frac{1}{2}\|h\|^2$

$g = f \cdot x$ $\qquad\qquad h = g - x$

$\dfrac{dL}{dh} = h$

$\dfrac{dh}{dg} = \dfrac{dL}{dh} \cdot \dfrac{dh}{dg} = h \cdot I = h$

$\dfrac{dL}{df} = \dfrac{dL}{dg} \cdot \dfrac{dg}{df} = h \cdot x^T \quad \left[ \in \mathbb{R}^{n \times n} \right]$
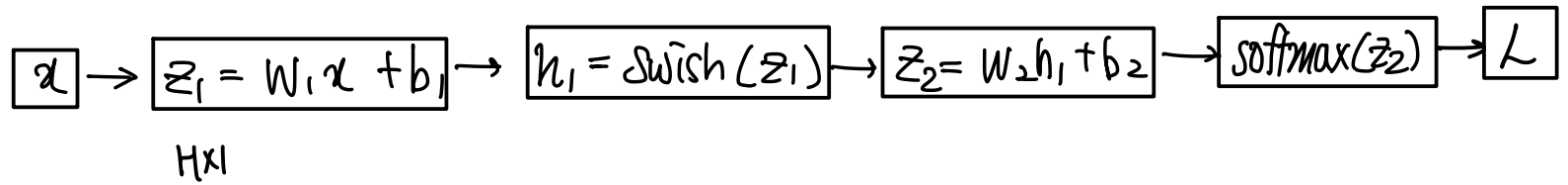
$\dfrac{dL}{dw} = \dfrac{df}{dw} \cdot \dfrac{dh}{df} = 2w \cdot (h \cdot x^T) \quad \left[ \in \mathbb{R}^{m \times n} \right]$

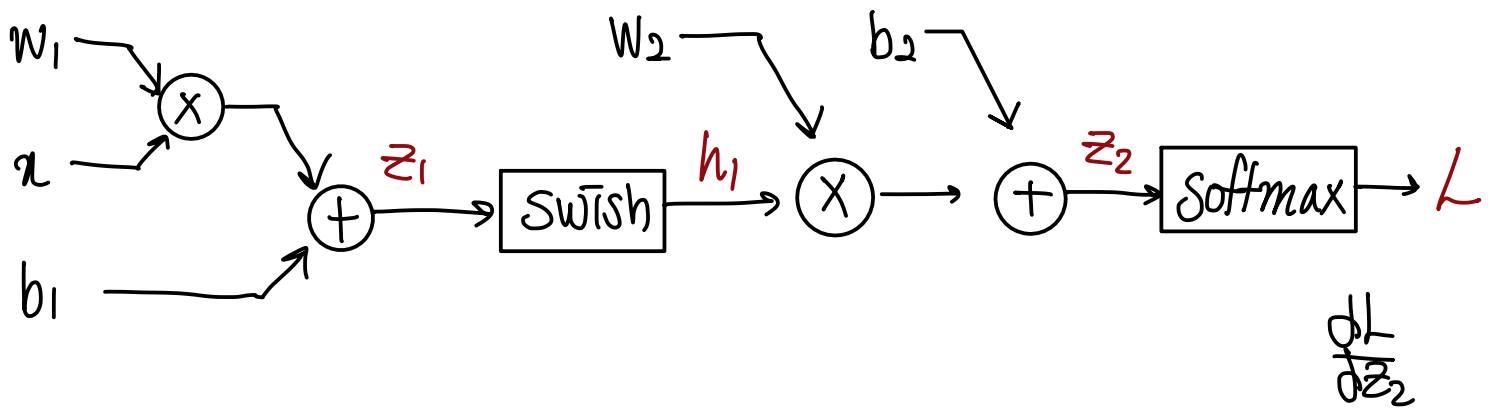$\Rightarrow \nabla_w L = 2w \cdot (W^T W x - x) x^T$

#2

I am a C147 Student.

#3.

$$x \rightarrow \boxed{z_1 = W_1 x + b_1} \rightarrow \boxed{h_1 = \text{Swish}(z_1)} \rightarrow \boxed{z_2 = W_2 h_1 + b_2} \rightarrow \boxed{\text{softmax}(z_2)} \rightarrow \boxed{L}$$

$H \times 1$

(a)

$$x \in \mathbb{R}^D, \ z_1 \in \mathbb{R}^H, \ W_1 \in \mathbb{R}^{H \times D}, \ b_1 \in \mathbb{R}^H, \ h_1 \in \mathbb{R}^H, \ z_2 \in \mathbb{R}^C, \ W_2 \in \mathbb{R}^{C \times H}, \ b_2 \in \mathbb{R}^C$$



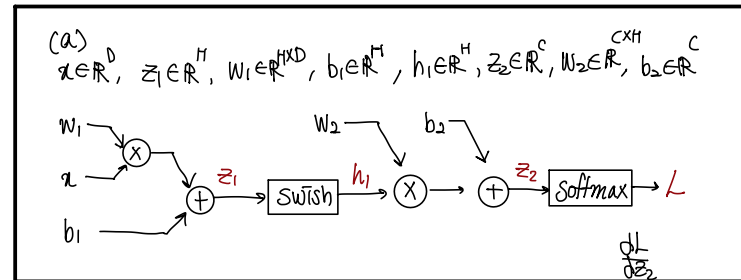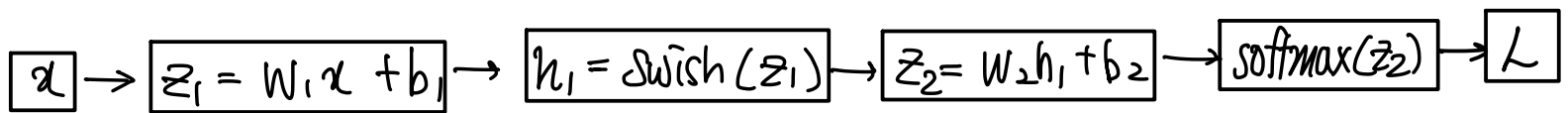(b)  $\nabla_{W_2} L, \ \nabla_{b_2} L$ ,,

$$z_2 = W_2 h_1 + b_2$$

$$\frac{dL}{db_2} = \frac{dL}{dz_2} \cdot \frac{dz_2}{db_2} = \frac{dL}{dz_2}$$

$$\frac{dL}{dW_2} = \frac{dL}{dz_2} \cdot \frac{dz_2}{dW_2} = \frac{dL}{dz_2} \cdot h_1^T \quad [\mathbb{R}^{C \times H}]$$

$$\Rightarrow \quad \nabla_{W_2} L = \frac{dL}{dz_2} \cdot h_1^T, \quad \nabla_{b_2} L = \frac{dL}{dz_2}$$

(c) $\nabla_{W_1} L$, $\nabla_{b_1} L$ "

$$\boxed{x} \rightarrow \boxed{z_1 = W_1 x + b_1} \rightarrow \boxed{h_1 = Swish(z_1)} \rightarrow \boxed{z_2 = W_2 h_1 + b_2} \rightarrow \boxed{softmax(z_2)} \rightarrow \boxed{L}$$

(a)
$x \in \mathbb{R}^D$, $z_1 \in \mathbb{R}^H$, $W_1 \in \mathbb{R}^{H \times D}$, $b_1 \in \mathbb{R}^H$, $h_1 \in \mathbb{R}^H$, $z_2 \in \mathbb{R}^C$, $W_2 \in \mathbb{R}^{C \times H}$, $b_2 \in \mathbb{R}^C$

$$W_1 \searrow \otimes$$
$$x \rightarrow$$
$$b_1 \nearrow$$
$$\oplus \xrightarrow{z_1} \boxed{Swish} \xrightarrow{h_1} \otimes \rightarrow \oplus \xrightarrow{z_2} \boxed{softmax} \rightarrow L$$
$$W_2 \qquad b_2$$
$$\frac{dL}{dz_2}$$

$$h_1 = z_1 \cdot \sigma(z_1)$$

$$\frac{dh_1}{dz_1} = \sigma(z_1) + z_1 \sigma(z_1)(1 - \sigma(z_1)) \quad \left[\in \mathbb{R}^{C \times H}\right]$$

$$\frac{dL}{dh_1} = \frac{dz_2}{dh_1} \cdot \frac{dL}{dz_2} = W_2^T \cdot \frac{dL}{dz_2} \quad \left[\in \mathbb{R}^H\right]$$

$$\frac{dL}{dz_1} = \frac{dh_1}{dz_1} \cdot \frac{dL}{dh_1} = \left[\sigma(z_1) + z_1 \sigma(z_1)(1 - \sigma(z_1))\right] \cdot \left(W_2^T \cdot \frac{dL}{dz_2}\right) \quad \left[\in \mathbb{R}^C\right]$$

$$\frac{dL}{dW_1} = \frac{dL}{dz_1} \cdot \frac{dz_1}{dW_1} = \frac{dL}{dz_1} \cdot x^T \quad \left[\in \mathbb{R}^{H \times D}\right]$$

$$\frac{dL}{db_1} = \frac{dL}{dz_1} \cdot \frac{dz_1}{db_1} = \frac{dL}{dz_1} \quad \left[\in \mathbb{R}^C\right]$$

$$\Rightarrow \nabla_{W_1} L = \left[\sigma(z_1) + z_1 \sigma(z_1)(1 - \sigma(z_1))\right] \cdot x^T$$

$$\nabla_{b_1} L = \left[\sigma(z_1) + z_1 \sigma(z_1)(1 - \sigma(z_1))\right]$$

, where $\sigma(k)$ is the sigmoid activation function.