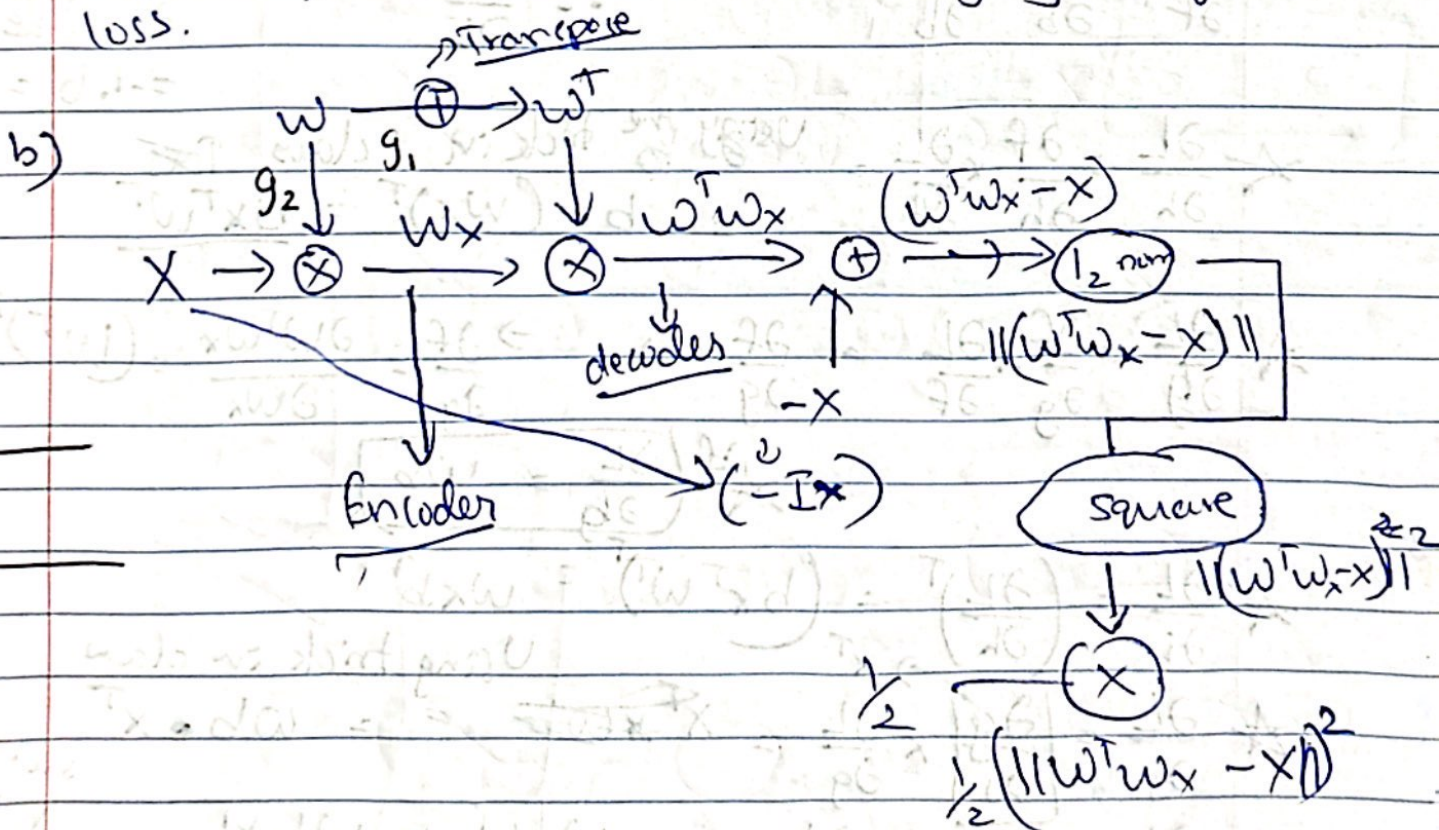


## Homework 3

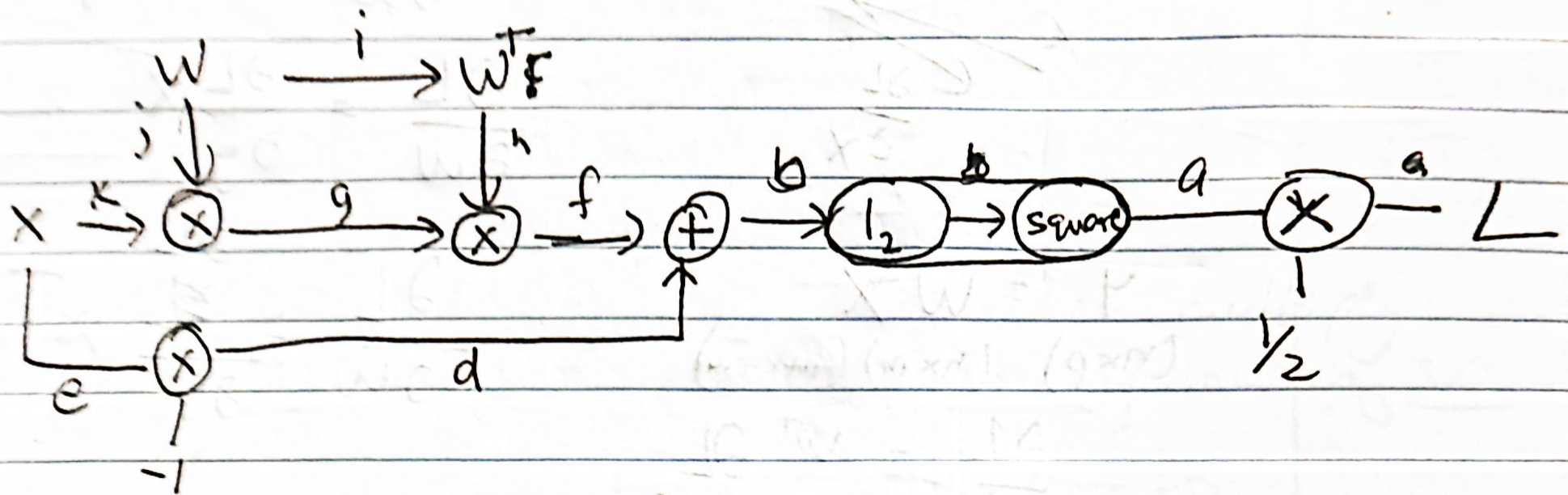
KRISH PATEL

(Q1)  $W \in \mathbb{R}^{m \times n}$ 

- a) When  $W$  is a square matrix, (assumption)  $W^T W$  would minimize the value of  $x$  the transformation (If  $W^T$  is orthogonal,  $W^T W$  would be equal to  $I$ , thus minimizing  $(W^T W x - x)$ ). Now, moving to this example,  $W$  reduces the dimensionality of  $x$ ,  $W^T$  is used for the reconstruction (similar to PCA). ~~By this~~ By minimizing the difference of the  $x'$  (reduced  $x$  dimensionality, which is the reconstructed to  $x$ 's original dimension) and  $x$ , we are minimizing information loss.



count





c) W has 2 paths, one through  $Wx$  and one through  $W \rightarrow \textcircled{1} \rightarrow W^T \rightarrow W^T Wx$ . These converge at  $W^T Wx$ . Thus, we can just add this  $\rightarrow$

$$\frac{\partial L}{\partial W} = \frac{\partial g_1}{\partial W} \cdot \frac{\partial L}{\partial g_1} + \frac{\partial g_2}{\partial W} \cdot \frac{\partial L}{\partial g_2} \quad \text{where } g_1 \text{ and } g_2 \text{ are defined in the figure behind.}$$

$$a = \frac{1}{2} \|b - \hat{b}\|^2 = \frac{1}{2} \|b - Wx\|^2$$

d)  $\frac{\partial L}{\partial a} = \frac{1}{2}$   $\frac{\partial L}{\partial b} = \frac{\partial a}{\partial b} \times \frac{\partial L}{\partial a} = 1 \times \frac{1}{2} = \frac{1}{2}$

$\star \frac{\partial L}{\partial f} = \frac{\partial f}{\partial b} \times \frac{\partial L}{\partial f} = 1 \times \frac{1}{2} = \frac{1}{2} \therefore \frac{\partial L}{\partial b} = \frac{1}{2}$   $\frac{\partial L}{\partial c} = \frac{\partial d}{\partial c} \times \frac{\partial L}{\partial d} = -1 \times \frac{1}{2} = -\frac{1}{2}$

$\star \frac{\partial L}{\partial h} = \frac{\partial f}{\partial h} \times \frac{\partial L}{\partial f} = \text{Using the trick in class, } \frac{\partial f}{\partial h} = b \times (Wx)^T = b x^T W^T$

$\star \frac{\partial L}{\partial g} = \frac{\partial f}{\partial g} \times \frac{\partial L}{\partial f} = \frac{\partial f}{\partial g} \times b \Rightarrow \frac{\partial f}{\partial g} = \frac{\partial W^T Wx}{\partial Wx} = (W^T)^T = W$

$\star \frac{\partial L}{\partial g} = Wb$

$\star \frac{\partial L}{\partial i} = \left( \frac{\partial L}{\partial h} \right)^T = (b x^T W^T)^T = Wx b^T$  Using trick in class

$\star \frac{\partial L}{\partial j} = \frac{\partial g}{\partial j} \times \frac{\partial L}{\partial g} = x^T W^T = Wb x^T$

$\star \frac{\partial L}{\partial W} = \frac{\partial L}{\partial j} + \frac{\partial L}{\partial i} = Wx b^T + Wb x^T$  substituting  $b = W^T Wx - x$

$$= Wx (W^T Wx - x)^T + W (W^T Wx - x) x^T$$



Q2) I am a ECE 147 student

Q3) NNDL to the Rescue!!!!

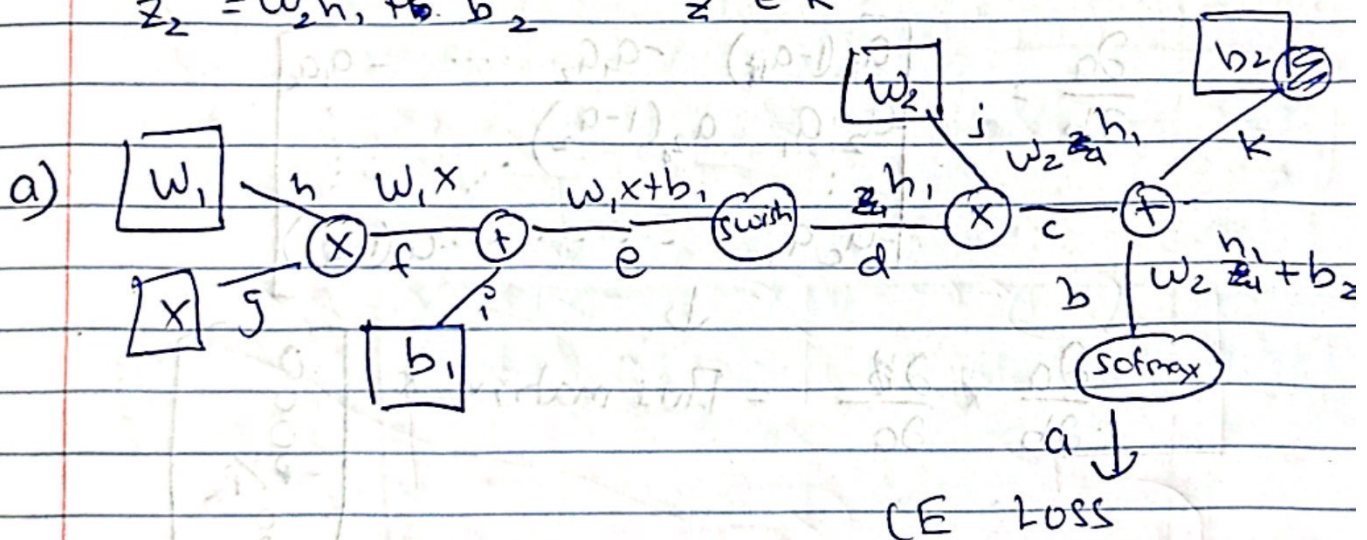
Input layer  $\rightarrow$   $D$  neurons.  
 Hidden layer  $\rightarrow$   $H$  neurons  
 Output layer  $\rightarrow$   $C$  neurons = 7

$$\text{Swish}(k) = \frac{k}{1+e^{-k}} = k \underbrace{\sigma(k)}_{\text{sigmoid}}$$

$$z_1 = w_1 x + b_1, \quad z_1 \in \mathbb{R}^H$$

$$h_1 = \text{swish}(z_1)$$

$$z_2 = w_2 h_1 + b_2, \quad z_2 \in \mathbb{R}^C$$



b)

$$\nabla_{w_2} L = \frac{\partial L}{\partial a}, \quad \nabla_{b_2} L = \frac{\partial L}{\partial b}$$

where  $b = w_2 h_1 + b_2$

$$\frac{\partial L}{\partial a} = \frac{\partial (-\log p_k)}{\partial a} = \frac{\partial L}{\partial b} \times \frac{\partial b}{\partial a} \quad \text{ON NEXT PAGE}$$

where  $\frac{\partial L}{\partial a}$  is for the correct class.

$$\frac{\partial L}{\partial a} = [0 \ 0 \ \dots \ -1/p_k \ 0 \ 0 \ 0]^T$$



$$\frac{\partial a}{\partial b} = \begin{bmatrix} \quad \quad \quad \end{bmatrix} \in \mathbb{R}^{C \times C}$$

↓  
Jacobian matrix

$$a_{ij} = \frac{e^{b_{ij}}}{\sum_{k=1}^C e^{b_{ik}}}$$

When  $j=1$   $\frac{\partial a_{ij}}{\partial b_{i1}} = a_{ij}(1-a_{ij})$

$j \neq 1$   $\frac{\partial a_{ij}}{\partial b_{i1}} = -a_{ij} \cdot a_{i1}$

$$\frac{\partial a}{\partial b} = \begin{bmatrix} a_{11}(1-a_{11}) & -a_{11}a_{12} & \dots & -a_{11}a_{1C} \\ a_{21}a_{11} & a_{21}(1-a_{21}) & \dots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ -a_{C1}a_{11} & \dots & \dots & a_{C1}(1-a_{C1}) \end{bmatrix}$$

↓

$$\frac{\partial a}{\partial b} \times \frac{\partial H}{\partial a} = \text{This matrix} \times \begin{bmatrix} 0 \\ 0 \\ 0 \\ -1/p_k \\ 0 \end{bmatrix}$$

$$= a_k(1-a_k) \begin{bmatrix} -a_{1k} \\ -a_{2k} \\ \vdots \\ a_k(1-a_k) \\ \vdots \\ -a_{Ck} \end{bmatrix} \times \frac{-1}{p_k}$$



$$= \begin{bmatrix} -\sigma(x_1)\sigma(x_k) \\ -\sigma(x_2)\sigma(x_k) \\ \vdots \\ \sigma(x_k)(1-\sigma(x_k)) \\ \vdots \\ \sigma(x_c)\sigma(x_k) \end{bmatrix} \begin{pmatrix} -1 \\ \sigma_k \end{pmatrix} = \begin{bmatrix} \sigma(x_1) \\ \sigma(x_2) \\ \vdots \\ 1(1-\sigma(x_k)) \\ \vdots \\ \sigma(x_c) \end{bmatrix}$$

$$= \text{softmax}(b) - \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

THUS, !!!

$$\frac{\partial L}{\partial b} = \text{softmax}(b) - e_k$$

Substituting  $b = z_2$

$$\frac{\partial L}{\partial b_2} = \frac{\partial b}{\partial b_2} \times \frac{\partial L}{\partial b} = 1 \times (\text{softmax}(z_2) - e_k) \times \text{forward pass derivative}$$

$\downarrow$   $\downarrow$   
 $z_2$   $\text{forward pass derivative}$

$$\frac{\partial L}{\partial b_2} = \begin{pmatrix} \frac{\partial L}{\partial z_j} \\ \downarrow \\ h_1 \end{pmatrix} \times \frac{\partial b}{\partial c} \times \frac{\partial L}{\partial b} = \text{softmax}(z_2) - e_k \times h_1^T$$

$\downarrow$   $\downarrow$   
 $1$  shortcut Using the weird rule in class.

c) Now that we have  $\frac{\partial L}{\partial d} = w_2^T (\text{softmax}(b) - e_k)$

we can back prop.  
Back prop through swish.  $\sigma$



$$\text{swish}(k) = \frac{k}{1+e^{-k}}$$

derivative:  $\frac{\partial (\text{swish}(k))}{\partial k}$

$$= 1 - \frac{k}{1+e^{-k}}$$

$$= \frac{1(1+e^{-k}) + k(e^{-k})}{(1+e^{-k})^2}$$

$$= \left( \frac{1}{\sigma(k)} + k \left( \frac{1}{\sigma(k)} - 1 \right) \right) \sigma(k)^2$$

$$= \sigma(k) + k(\sigma(k) - \sigma(k)^2)$$

$$= \sigma(k) +$$

Substituting  $z_1 = \underline{w_1 x + b}$  and  $h_1 = \text{swish}(w_1 x + b)$  and  $k = z_1$

$$\sigma(k) + k \sigma(k) (1 - \sigma(k))$$

$$= \left( h_1 + \sigma(z_1) \odot (1 - h_1) \right) \left[ \begin{aligned} &\sigma(z_1) + h_1 - h_1 \times \sigma(z_1) \\ &= (h_1 + \sigma(z_1)(1 - h_1)) \end{aligned} \right]$$

Now,  $\frac{\partial L}{\partial h} = \frac{\partial f}{\partial h} \times \frac{\partial L}{\partial f} =$

$$= \left( (h_1 + \sigma(z_1) \odot (1 - h_1)) \odot (W_2^T (\text{softmax}(z_2) - e_k)) \right)^T x$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial f} \cdot \left( (h_1 + \sigma(z_1) \odot (1 - h_1)) \odot (W_2^T (\text{softmax}(z_2) - e_k)) \right)$$