

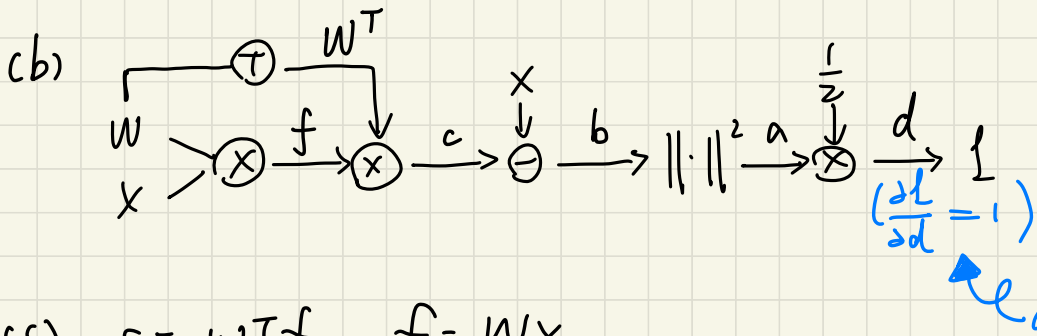
1. //

(a) If we use  $Wx$  in the loss function, the dimension of  $Wx$  is  $(m \times n) \times (n \times 1) = (m \times 1)$ , where  $m < n$ .

This means some information about  $x$  is lost.

However, if we use  $W^T W x$  here, the dimension of  $W^T W x$  is  $(n \times m) \times (m \times n) \times (n \times 1) = (n \times 1)$ , which is the same as  $\dim(x)$ .

In this way, the minimization can find a  $W$  ought to preserve information about  $x$ .



Computing the local gradient,

$$\frac{\partial c}{\partial f} = (W^T)^T = W, \quad \frac{\partial l}{\partial f} = W \frac{\partial l}{\partial c}.$$

$$\left(\frac{\partial l}{\partial W}\right)_1 = \frac{\partial f}{\partial W} \cdot \frac{\partial l}{\partial f} = \frac{\partial l}{\partial f} \cdot x^T = W \frac{\partial l}{\partial c} x^T.$$

Now for the second path  $W^T$ :

$$\left(\frac{\partial l}{\partial W}\right)_2 = \left(\frac{\partial l}{\partial W^T}\right)^T = \left(\frac{\partial c}{\partial W^T} \cdot \frac{\partial l}{\partial c}\right)^T = \left(\frac{\partial l}{\partial c} \cdot W^T x^T\right)^T = Wx \left(\frac{\partial l}{\partial c}\right)^T.$$

$$\text{So, } \nabla_W L = \left(\frac{\partial l}{\partial W}\right)_1 + \left(\frac{\partial l}{\partial W}\right)_2 = W \frac{\partial l}{\partial c} x^T + Wx \left(\frac{\partial l}{\partial c}\right)^T.$$

(d) since  $\otimes$  switches the gradient,

$$\frac{\partial l}{\partial a} = \frac{1}{2} \frac{\partial l}{\partial a} = \frac{1}{2}.$$

Now,  $a = \|b\|_2^2$

$$\frac{\partial a}{\partial b} = 2b \in \mathbb{R}^n, \quad \frac{\partial l}{\partial b} = b \in \mathbb{R}^n$$

since  $\ominus$  passes the gradient,

$$\frac{\partial L}{\partial c} = \frac{\partial L}{\partial b} = b \in \mathbb{R}^n$$

using the result from question (c),

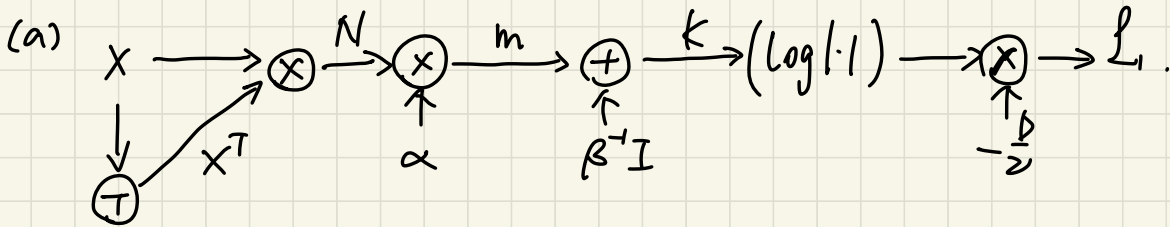
$$\begin{aligned} \tau_W L &= W \frac{\partial L}{\partial c} X^T + W X \left( \frac{\partial L}{\partial c} \right)^T \\ &= W \cdot b \cdot X^T + W X \cdot b^T \end{aligned}$$

$$b = W^T W X - X,$$

$$\text{so } \tau_W L = W(W^T W X - X) X^T + W X (W^T W X - X)^T$$

$$\text{2. } L_1 = -\frac{D}{2} \log |\alpha X X^T + \beta^{-1} I|$$

$$L_2 = -\frac{1}{2} \text{tr}((\alpha X X^T + \beta^{-1} I)^{-1} Y Y^T).$$



(b)  $L_1 = -\frac{D}{2} \log(\det(K))$

using equation 57 from cookbook,

$$\frac{\partial L_1}{\partial K} = -\frac{D}{2} (K^T)^{-1} = -\frac{D}{2} (K^T)^{-1}$$

since  $\oplus$  passes the gradient,

$$\frac{\partial L_1}{\partial m} = \frac{\partial L_1}{\partial K} = -\frac{D}{2} (K^T)^{-1}$$

since  $\otimes$  switches the gradient,

$$\frac{\partial L_1}{\partial N} = -\frac{D}{2} (K^T)^{-1}$$

$N = X X^T$ , using the hint,

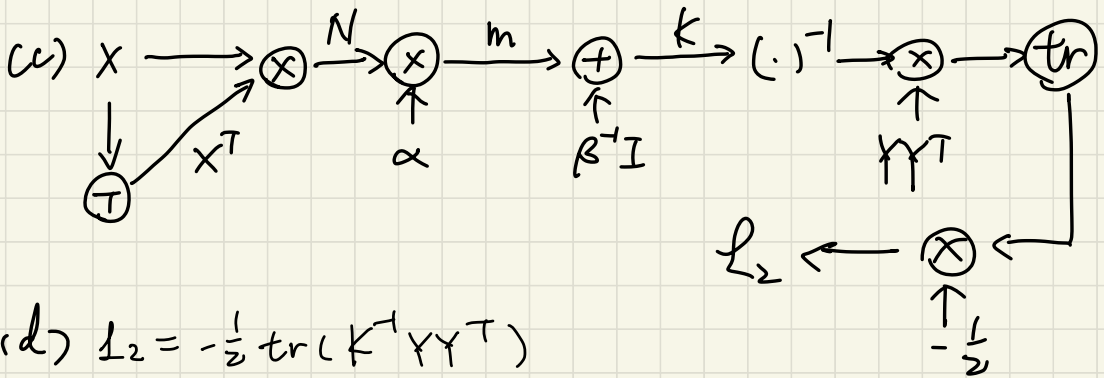
$$\left(\frac{\partial L_1}{\partial X}\right)_1 = \frac{\partial N}{\partial X} \cdot \frac{\partial L_1}{\partial N} = \frac{\partial L_1}{\partial N} (X^T)^T = -\frac{D}{2} (K^T)^{-1} X.$$

$$\left(\frac{\partial L_1}{\partial X^T}\right)_2 = -\frac{D}{2} (K^T X)^T$$

$$\text{so } \frac{\partial L_1}{\partial X} = \left(\frac{\partial L_1}{\partial X}\right)_1 + \left(\frac{\partial L_1}{\partial X^T}\right)_2^T = -\frac{D}{2} (K^T)^{-1} X - \frac{D}{2} K^T X$$

since  $K = \alpha XX^T + \beta^T I$ ,  $K$  is symmetric.

$$\frac{\partial \mathcal{L}_1}{\partial X} = -\alpha D K^T X \quad (\text{where } K = \alpha XX^T + \beta^T I)$$



(d)  $\mathcal{L}_2 = -\frac{1}{2} \text{tr}(K^T YY^T)$

Using Equation 124 from cookbook,

$$\frac{\partial \mathcal{L}_2}{\partial K} = \frac{1}{2} (K^T YY^T K^T)^T$$

Since  $\oplus$  passes the gradient,

$$\frac{\partial \mathcal{L}_2}{\partial m} = \frac{\partial \mathcal{L}_2}{\partial K}$$

since  $\otimes$  switches the gradient,

$$\frac{\partial \mathcal{L}_2}{\partial N} = \frac{\alpha}{2} (K^T YY^T K^T)^T$$

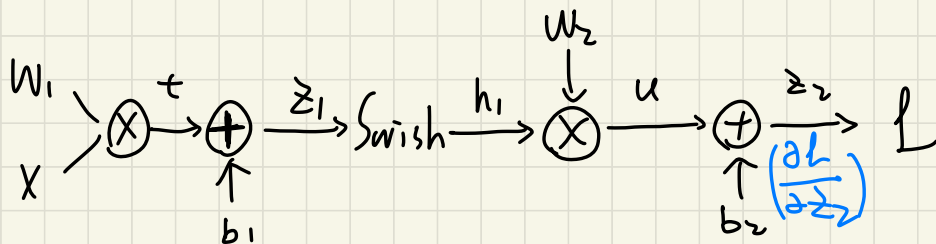
$$\begin{aligned} \frac{\partial \mathcal{L}_2}{\partial X} &= \left( \frac{\partial \mathcal{L}_2}{\partial X} \right)_1 + \left( \frac{\partial \mathcal{L}_2}{\partial X^T} \right)_2^T = \frac{\partial \mathcal{L}_2}{\partial N} \cdot X + \left( X^T \frac{\partial \mathcal{L}_2}{\partial N} \right)^T \\ &= \frac{\alpha}{2} (K^T YY^T K^T)^T X + \frac{\alpha}{2} (K^T YY^T K^T) X \\ &= \alpha K^T YY^T K^T X. \end{aligned}$$

$$\frac{\partial \mathcal{L}_2}{\partial X} = \alpha K^T YY^T K^T X \quad (\text{where } K = \alpha XX^T + \beta^T I)$$

(e)  $\frac{\partial \mathcal{L}}{\partial X} = \frac{\partial \mathcal{L}_1}{\partial X} + \frac{\partial \mathcal{L}_2}{\partial X} = \alpha (K^T YY^T K^T X - D K^T X)$

3.

(a)



(b) since  $\oplus$  passes the gradient,

$$\nabla_{b_2} L = \frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial z_2} \in \mathbb{R}^c$$

$$\frac{\partial L}{\partial u} = \frac{\partial L}{\partial z_2}$$

Now,  $u = w_2 h_1$ .

$$\nabla_{w_2} L = \frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial u} \cdot h_1^T = \frac{\partial L}{\partial z_2} \cdot h_1^T \in \mathbb{R}^{c \times H}$$

$$(c) \frac{\partial L}{\partial h_1} = \frac{\partial u}{\partial h_1} \cdot \frac{\partial L}{\partial u} = w_2^T \cdot \frac{\partial L}{\partial u} = w_2^T \frac{\partial L}{\partial z_2}$$

Now,  $h_1 = \text{swish}(z_1) = z_1 \sigma(z_1)$

$$\begin{aligned} \frac{\partial h_1}{\partial z_1} &= z_1 \frac{\partial \sigma(z_1)}{\partial z_1} + \sigma(z_1) \\ &= z_1 \sigma(z_1) [1 - \sigma(z_1)] + \sigma(z_1) \\ &= \sigma(z_1) [z_1 + 1 - z_1 \sigma(z_1)] \end{aligned}$$

$$\frac{\partial L}{\partial z_1} = \frac{\partial h_1}{\partial z_1} \cdot \frac{\partial L}{\partial h_1}$$

Since  $\oplus$  passes the gradient,

$$\nabla_{b_1} L = \frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial z_1} = \sigma(z_1) [z_{1+1} - z_1 \sigma(z_1)] w_2^T \frac{\partial L}{\partial z_2}$$

$$\frac{\partial L}{\partial t} = \frac{\partial L}{\partial z_1}$$

$$\text{Now } t = w_1 x$$

$$\nabla_{w_1} L = \frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial t} \cdot x^T = \sigma(z_1) [z_{1+1} - z_1 \sigma(z_1)] w_2^T \frac{\partial L}{\partial z_2} \cdot x^T$$