

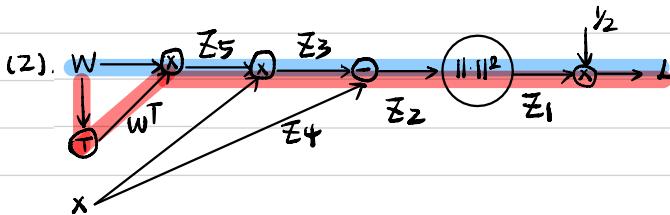
1. (15 points) **Backpropagation for autoencoders.** In an autoencoder, we seek to reconstruct the original data after some operation that reduces the data's dimensionality. We may be interested in reducing the data's dimensionality to gain a more compact representation of the data.

For example, consider $\mathbf{x} \in \mathbb{R}^n$. Further, consider $\mathbf{W} \in \mathbb{R}^{m \times n}$ where $m < n$. Then \mathbf{Wx} is of lower dimensionality than \mathbf{x} . One way to design \mathbf{W} so that \mathbf{Wx} still contains key features of \mathbf{x} is to minimize the following expression

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

with respect to \mathbf{W} . (To be complete, autoencoders also have a nonlinearity in each layer, i.e., the loss is $\frac{1}{2} \|f(\mathbf{W}^T f(\mathbf{Wx})) - \mathbf{x}\|^2$. However, we'll work with the linear example.)

(1) Because if $L \rightarrow 0$, then it means $\mathbf{W}^T \mathbf{Wx} - \mathbf{x} \rightarrow 0$, regard $\mathbf{W}^T \mathbf{W}$ as a function $f(\cdot)$. Then $f(\mathbf{x}) \rightarrow \mathbf{x}$, then it means after transformation, $f(\mathbf{x})$ is still very close to original \mathbf{x} , it means information of \mathbf{x} has been preserved.



(2). There are \mathbf{w} and \mathbf{w}^T two paths to account for $\frac{\partial L}{\partial \mathbf{w}}$.

$$\frac{\partial L}{\partial \mathbf{w}} = (\frac{\partial L}{\partial \mathbf{w}})_{\text{blue}} + (\frac{\partial L}{\partial \mathbf{w}^T})_{\text{pink}}^T$$

$$(3). L = \frac{1}{2} z_1, \frac{\partial L}{\partial z_1} = \frac{1}{2}, z_1 = \|z_2\|^2, \frac{\partial z_1}{\partial z_2} = 2z_2, \frac{\partial L}{\partial z_2} = \frac{\partial z_1}{\partial z_2} \cdot \frac{\partial L}{\partial z_1} = z_2$$

"Add" passes gradient: $\frac{\partial L}{\partial z_3} = \frac{\partial L}{\partial z_2}$, $z_3 = z_5 \cdot x$, $\frac{\partial L}{\partial z_5} = \frac{\partial z_3}{\partial z_5} \cdot \frac{\partial L}{\partial z_3} = \frac{\partial L}{\partial z_3} \cdot x^T$

$$(\frac{\partial L}{\partial \mathbf{w}})_1 = \frac{\partial z_5}{\partial \mathbf{w}} \cdot \frac{\partial L}{\partial z_5} = \mathbf{w} \cdot \frac{\partial L}{\partial z_5} = \mathbf{w} \cdot z_2 \cdot x^T, (\frac{\partial L}{\partial \mathbf{w}^T})_2 = \frac{\partial z_5}{\partial \mathbf{w}^T} \frac{\partial L}{\partial z_5} = \frac{\partial L}{\partial z_5} \cdot \mathbf{w}^T = z_2 \cdot x^T \cdot \mathbf{w}^T$$

$$z_5 = \mathbf{w}^T \mathbf{w}$$

$$\Rightarrow \frac{\partial L}{\partial \mathbf{w}} = (\frac{\partial L}{\partial \mathbf{w}})_1 + (\frac{\partial L}{\partial \mathbf{w}^T})_2^T = \mathbf{w} \cdot z_2 \cdot x^T + \mathbf{w} \cdot x \cdot z_2^T = \mathbf{w} \left((\mathbf{w}^T \mathbf{w} \mathbf{x} - \mathbf{x}) \mathbf{x}^T + \mathbf{x} (\mathbf{w}^T \mathbf{w} \mathbf{x} - \mathbf{x})^T \right)$$

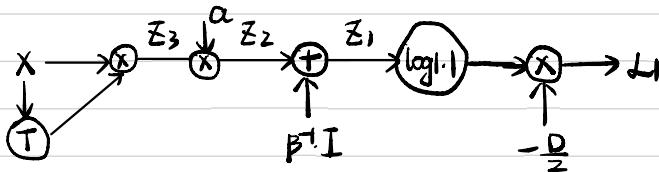
$$z_2 = \mathbf{w}^T \mathbf{w} \mathbf{x} - \mathbf{x}$$

2. $\mathcal{L} = -c - \frac{D}{2} \log |\mathbf{K}| - \frac{1}{2} \text{tr}(\mathbf{K}^{-1} \mathbf{Y} \mathbf{Y}^T)$ $\mathbf{K} = \alpha \mathbf{X} \mathbf{X}^T + \beta^{-1} \mathbf{I}$

$$\mathcal{L}_1 = -\frac{D}{2} \log |\alpha \mathbf{X} \mathbf{X}^T + \beta^{-1} \mathbf{I}|$$

$$\mathcal{L}_2 = -\frac{1}{2} \text{tr}((\alpha \mathbf{X} \mathbf{X}^T + \beta^{-1} \mathbf{I})^{-1} \mathbf{Y} \mathbf{Y}^T)$$

(1). Draw Computation Graph for \mathcal{L}_1 .



(2) Compute $\frac{\partial \mathcal{L}_1}{\partial \mathbf{x}}$: $\mathcal{L}_1 = -\frac{D}{2} \log |\det \mathbf{Z}_1|$, $\frac{\partial \mathcal{L}_1}{\partial \mathbf{z}_1} = -\frac{D}{2} (\mathbf{Z}_1^T)^{-1} = -\frac{D}{2} (\mathbf{Z}_1^T)^{-1}$.
 $\frac{\partial \mathcal{L}_1}{\partial \mathbf{z}_2} = \frac{\partial \mathcal{L}_1}{\partial \mathbf{z}_1}$, $\frac{\partial \mathcal{L}_1}{\partial \mathbf{z}_3} = \alpha \cdot \frac{\partial \mathcal{L}_1}{\partial \mathbf{z}_1} = -\frac{D\alpha}{2} (\mathbf{Z}_1^T)^{-1}$

$$\mathbf{Z}_3 = \mathbf{X} \mathbf{X}^T, \frac{\partial \mathcal{L}_1}{\partial \mathbf{X}} = \frac{\partial \mathbf{Z}_3}{\partial \mathbf{X}}, \frac{\partial \mathcal{L}_1}{\partial \mathbf{z}_3} = \frac{\partial \mathcal{L}_1}{\partial \mathbf{z}_3} \cdot (\mathbf{X}^T)^T = \frac{\partial \mathcal{L}_1}{\partial \mathbf{z}_3} \cdot \mathbf{X}$$

$$\frac{\partial \mathcal{L}_1}{\partial \mathbf{X}^T} = \frac{\partial \mathbf{Z}_3}{\partial \mathbf{X}^T}, \frac{\partial \mathcal{L}_1}{\partial \mathbf{z}_3} = \mathbf{X}^T \cdot \frac{\partial \mathcal{L}_1}{\partial \mathbf{z}_3}$$

$$\Rightarrow \frac{\partial \mathcal{L}_1}{\partial \mathbf{X}} = \left(\frac{\partial \mathcal{L}_1}{\partial \mathbf{X}} + \left(\frac{\partial \mathcal{L}_1}{\partial \mathbf{X}^T} \right)^T \right) \mathbf{X} = \frac{\partial \mathcal{L}_1}{\partial \mathbf{z}_3} \cdot \mathbf{X} + \left(\frac{\partial \mathcal{L}_1}{\partial \mathbf{z}_3} \right)^T \mathbf{X}$$

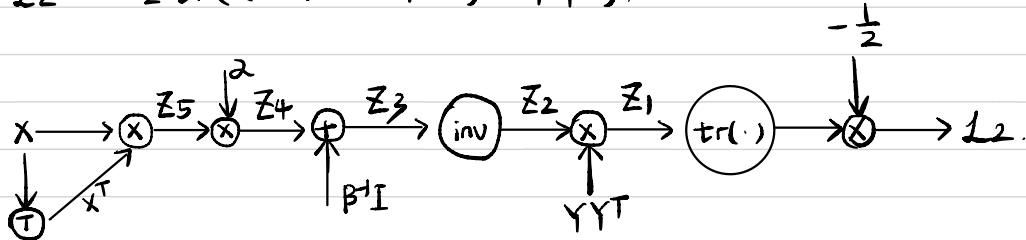
$$= -\frac{D\alpha}{2} (\mathbf{Z}_1^T)^{-1} \cdot \mathbf{X} + \left(-\frac{D\alpha}{2} \cdot (\mathbf{Z}_1^T)^T \right)^T \mathbf{X}$$

$\mathbf{Z}_1 = \mathbf{X} \mathbf{X}^T \alpha + \beta^{-1} \mathbf{I}$ is symmetric, $\frac{\partial \mathcal{L}_1}{\partial \mathbf{X}} = -\frac{D\alpha}{2} (\mathbf{Z}_1^T)^{-1} \mathbf{X} - \frac{D\alpha}{2} \cdot \mathbf{Z}_1^{-1} \mathbf{X} = -D\alpha \mathbf{Z}_1^{-1} \mathbf{X}$

$$= -D\alpha \cdot (\mathbf{X} \mathbf{X}^T \alpha + \beta^{-1} \mathbf{I})^{-1} \mathbf{X}$$

(3). Draw Computational Graph for L_2 .

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



(4) Compute $\frac{\partial L_2}{\partial X}$.

$$L_2 = -\frac{1}{2} \text{tr}(Z_1), \quad \frac{\partial L_2}{\partial Z_1} = -\frac{1}{2} - \frac{(Z_1)}{\partial Z_1} = -\frac{1}{2} I$$

$$Z_1 = Z_2(YY^T), \quad \frac{\partial L_2}{\partial Z_2} = \frac{\partial Z_2}{\partial Z_2} \cdot \frac{\partial L_2}{\partial Z_1} = -\frac{1}{2} I \cdot (YY^T)^T = -\frac{1}{2} YY^T.$$

$$Z_2 = (Z_3)^{-1}, \quad \frac{\partial Z_2}{\partial Z_3} = -Z_3^{-T} \cdot \frac{\partial Z_2}{\partial Z_3^{-1}} \cdot Z_3^{-T} = -Z_3^{-T} \cdot I \cdot Z_3^{-T} = -(Z_3^{-T})^2.$$

$$\frac{\partial L}{\partial Z_3} = \frac{\partial Z_2}{\partial Z_3} \cdot \frac{\partial L}{\partial Z_2} = + (Z_3^{-T})^2 \cdot \frac{1}{2} YY^T$$

"Add" passes gradient: $\frac{\partial L_2}{\partial Z_4} = \frac{\partial L_2}{\partial Z_3}$

$$Z_4 = \alpha \cdot Z_5, \quad \frac{\partial L_2}{\partial Z_5} = \alpha \cdot \frac{\partial L_2}{\partial Z_4} = \alpha \cdot \frac{\partial L_2}{\partial Z_3} = \frac{\alpha}{2} \cdot (Z_3^{-T} \cdot Z_3^{-T}) \cdot YY^T$$

$$Z_5 = XX^T, \quad \frac{\partial L_2}{\partial X} = \frac{\partial Z_5}{\partial X} \cdot \frac{\partial L_2}{\partial Z_5} = \frac{\partial L_2}{\partial Z_5} \cdot (X^T)^T = \frac{\partial L_2}{\partial Z_5} \cdot X$$

$$\frac{\partial L_2}{\partial X^T} = \frac{\partial Z_5}{\partial X^T} \cdot \frac{\partial L_2}{\partial Z_5} = X^T \cdot \frac{\partial L_2}{\partial Z_5}$$

$$\frac{\partial L_2}{\partial X} = (\frac{\partial L_2}{\partial X})_1 + (\frac{\partial L_2}{\partial X^T})^T = \frac{\partial L_2}{\partial Z_5} \cdot X + (\frac{\partial L_2}{\partial Z_5})^T \cdot X$$

$$= \frac{\alpha}{2} \cdot (Z_3^{-T} \cdot Z_3^{-T}) \cdot YY^T \cdot X + (\frac{\alpha}{2} \cdot (Z_3^{-T})^2 \cdot YY^T)^T \cdot X$$

$$= \frac{\alpha}{2} ((Z_3^{-T})^2 YY^T + YY^T \cdot (Z_3^{-T})^2) X$$

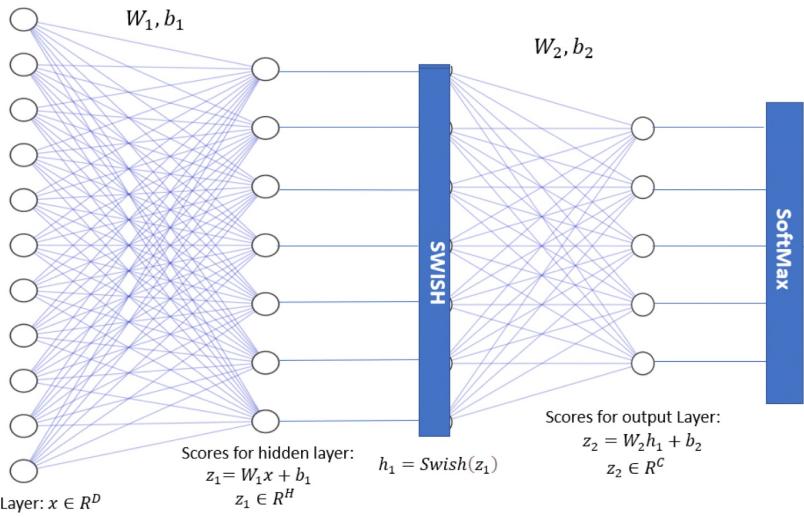
where $Z_3 = \alpha XX^T + \beta^T I$

$$(e). \frac{\partial L}{\partial X} = \frac{\partial}{\partial X} (-C + L_1 + L_2) = \frac{\partial L_1}{\partial X} + \frac{\partial L_2}{\partial X}$$

$$= -D \cdot a \cdot (XX^T \cdot a + \beta^T I)^{-1} X$$

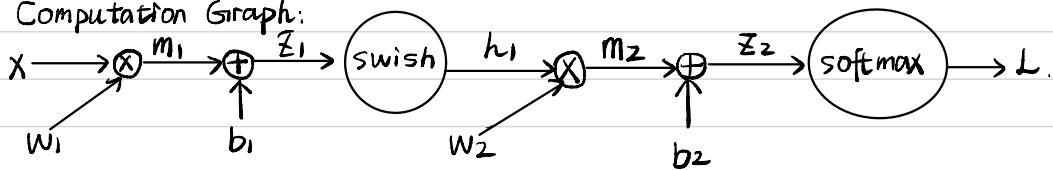
$$+ \frac{\alpha}{2} ((Z_3^{-T})^2 YY^T + YY^T \cdot (Z_3^{-T})^2) X$$

$$Z_3 = \alpha XX^T + \beta^T I.$$



$$\text{Swish}(k) = k \cdot \sigma(k)$$

(a). Computation Graph:



(b). m_1, m_2 defined as above.

"Add" passes gradient: $\frac{\partial L}{\partial m_2} = \frac{\partial L}{\partial z_2}, \frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial z_2}$.

$$m_2 = W_2 h_1, \frac{\partial L}{\partial W_2} = \frac{\partial m_2}{\partial W_2}, \frac{\partial L}{\partial m_2} = \frac{\partial L}{\partial m_2} \cdot h_1^T = \frac{\partial L}{\partial z_2} \cdot h_1^T$$

$$(c) \frac{\partial L}{\partial h_1} = \frac{\partial m_2}{\partial h_1} \frac{\partial L}{\partial m_2} = W_2^T \frac{\partial L}{\partial z_2}$$

$$h_1 = \text{swish}(z_1), \frac{\partial h_1}{\partial z_1} = \sigma(z_1) + z_1(1 - \sigma(z_1))$$

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

$$\frac{\partial L}{\partial z_1} = \frac{\partial h_1}{\partial z_1} \cdot \frac{\partial L}{\partial h_1} = \text{diag}(\sigma(z_1) + z_1(1 - \sigma(z_1))) \frac{\partial L}{\partial h_1} = \text{diag}(\sigma(z_1) + z_1(1 - \sigma(z_1))) \cdot W_2^T \frac{\partial L}{\partial z_2}$$

$$\text{"Add" passes gradient: } \frac{\partial}{\partial b_1} = \frac{\partial L}{\partial z_1}, \frac{\partial L}{\partial m_1} = \frac{\partial L}{\partial z_1}$$

$$m_1 = W_1 x, \frac{\partial L}{\partial W_1} = \frac{\partial m_1}{\partial W_1} \frac{\partial L}{\partial m_1} = \frac{\partial L}{\partial m_1} \cdot x^T = \frac{\partial L}{\partial z_1} \cdot x^T$$