

2. In Pytorch,

Weigthes $(A_L)_{ij}$ are initialized with $U\left(-\frac{1}{\sqrt{\text{fan_in}}}, \frac{1}{\sqrt{\text{fan_in}}}\right)$ and
biases $(b_L)_i$ are initialized with $U\left(-\frac{1}{\sqrt{\text{fan_in}}}, \frac{1}{\sqrt{\text{fan_in}}}\right)$.
They are mutually independently initialized.

$$E((A_L)_{ij}) = \frac{1}{2} \left(-\sqrt{\frac{1}{\text{fan_in}}} + \sqrt{\frac{1}{\text{fan_in}}} \right) = 0$$

$$\text{Var}((A_L)_{ij}) = \frac{1}{12} \left(\sqrt{\frac{1}{\text{fan_in}}} + \sqrt{\frac{1}{\text{fan_in}}} \right)^2 = \frac{1}{3} \cdot \frac{1}{\text{fan_in}} = \frac{1}{3} \cdot \frac{1}{n_{L-1}}$$

$$E((b_L)_i) = \frac{1}{2} \left(-\frac{1}{\sqrt{\text{fan_in}}} + \frac{1}{\sqrt{\text{fan_in}}} \right) = 0$$

$$\text{Var}((b_L)_i) = \frac{1}{12} \left(\frac{1}{\sqrt{\text{fan_in}}} + \frac{1}{\sqrt{\text{fan_in}}} \right)^2 = \frac{1}{3} \cdot \frac{1}{\text{fan_in}} = \frac{1}{3} \cdot \frac{1}{n_{L-1}}$$

$$\begin{aligned} \text{Var}(y_L) &= \text{Var}(A_L y_{L-1} + b_L) = \text{Var} \left(\sum_j \{(A_L)_{ij} (y_{L-1})_j\} + (b_L)_i \right) \\ &= n_{L-1} \text{Var}(A_L)_{ii} \text{Var}(y_{L-1})_i + \text{Var}(b_L)_i \\ &= n_{L-1} \frac{1}{3} \cdot \frac{1}{n_{L-1}} \text{Var}(y_{L-1})_i + \frac{1}{3} \cdot \frac{1}{n_{L-1}} \\ &= \frac{1}{3} \text{Var}(y_{L-1})_i + \frac{1}{3} \cdot \frac{1}{n_{L-1}} \end{aligned}$$

$$\begin{aligned} \text{Var}(y_L)_i &= \text{Var} \left(\sum_j \{(A_L)_{ij} (y_{L-1})_j\} + (b_L)_i \right) \\ &= n_{L-1} \text{Var}(A_L)_{ii} \text{Var}(y_{L-1})_i + \text{Var}(b_L)_i \\ &= n_{L-1} \frac{1}{3} \cdot \frac{1}{n_{L-1}} \text{Var}(y_{L-1})_i + \frac{1}{3} \cdot \frac{1}{n_{L-1}} \\ &= \frac{1}{3} \text{Var}(y_{L-1})_i + \frac{1}{3} \cdot \frac{1}{n_{L-1}}, \quad i = 1, 2, 3, \dots, L-1, \quad y_0 := x \end{aligned}$$

$$\therefore \text{Var}(y_L) = \frac{1}{3^L} \text{Var}(y_0) + \sum_{k=0}^{L-1} \frac{1}{3^{L-k}} \cdot \frac{1}{n_k}$$

$$= \frac{1}{3^L} + \sum_{k=0}^{L-1} \frac{1}{3^{L-k}} \cdot \frac{1}{n_k}$$

□

3.

$$\frac{\partial y_L}{\partial w_l} = \begin{bmatrix} (y_{l-1})_1 & (y_{l-1})_2 & \cdots & (y_{l-1})_{f_l} & 0 & \cdots & 0 \\ 0 & (y_{l-1})_2 & (y_{l-1})_3 & \cdots & (y_{l-1})_{f_l+1} & 0 & \cdots & 0 \\ \vdots & & & & & & & \\ 0 & \cdots & 0 & (y_{l-1})_{n_l} & (y_{l-1})_{n_l+1} & \cdots & (y_{l-1})_{f_l+n_l-1} \\ & & & & & & \end{bmatrix}_{n_{l-1}}$$

$$y'_l := A_l y_{l-1} + b_l, \quad l=1, \dots, L-1.$$

$$\frac{\partial y_L}{\partial w_l} = \frac{\partial y_L}{\partial y'_l} \frac{\partial y'_l}{\partial w_l} = \frac{\partial y_L}{\partial y'_l} \frac{\partial y'_l}{\partial y'_L} \frac{\partial y'_L}{\partial w_l}, \quad l=1, \dots, L-1$$

$$= \frac{\partial y_L}{\partial y'_l} \text{diag}(\sigma'(y'_l)) \begin{bmatrix} (y_{l-1})_1 & (y_{l-1})_2 & \cdots & (y_{l-1})_{f_l} \\ (y_{l-1})_2 & (y_{l-1})_3 & \cdots & (y_{l-1})_{f_l+1} \\ \vdots & & & \\ (y_{l-1})_{n_l} & \cdots & (y_{l-1})_{f_l+n_l-1} \\ & & \end{bmatrix}_{n_{l-1}}$$

$$= A_L \left[\text{diag}(\sigma'(y'_{L-1})) A_{L-1} \right] \cdots \left[\text{diag}(\sigma'(y'_{l+1})) A_{l+1} \right] \text{diag}(\sigma'(y'_l)) \begin{bmatrix} (y_{l-1})_1 & (y_{l-1})_2 & \cdots & (y_{l-1})_{f_l} \\ (y_{l-1})_2 & (y_{l-1})_3 & \cdots & (y_{l-1})_{f_l+1} \\ \vdots & & & \\ (y_{l-1})_{n_l} & \cdots & (y_{l-1})_{f_l+n_l-1} \\ & & \end{bmatrix}_{n_{l-1}}$$

$$\frac{\partial y_L}{\partial b_l} = I_{n_l}$$

Using result of Homework 4 problem 6,

$$\begin{aligned} \frac{\partial y_L}{\partial b_l} &= \frac{\partial y_L}{\partial y'_l} \frac{\partial y'_l}{\partial b_l} = \frac{\partial y_L}{\partial y_{l-1}} \frac{\partial y_{l-1}}{\partial y_{l-2}} \cdots \frac{\partial y_{l-1}}{\partial y'_l} \frac{\partial y'_l}{\partial b_l} \frac{\partial b_l}{\partial b_l} \\ &= A_L \left[\text{diag}(\sigma'(y'_{L-1})) A_{L-1} \right] \cdots \left[\text{diag}(\sigma'(y'_{l+1})) A_{l+1} \right] \text{diag}(\sigma'(y'_l)) \cdot I_{n_l} \\ &\quad l=1, \dots, L-1 \end{aligned}$$

Each $\text{diag} \cdot A$ is computed by $(A^T \cdot \text{diag})^T$ with transposed convolution.

By storing output of each layer during forward propagation,
back propagation can be done efficiently.

□

$$4. \quad Y = \varphi(x) = A^{-1}(x - b)$$

$$\begin{aligned}
P_X(\alpha) &= P_Y(\varphi(\alpha)) \left| \det \frac{\partial \varphi}{\partial \alpha}(\alpha) \right| \\
&= P_Y(A^{-1}(x - b)) | \det A^{-1} | \\
&= \frac{1}{(2\pi)^n} e^{-\frac{1}{2} (x-b)^T (A^{-1})^T A^{-1} (x-b)} |A|^{-1}, \quad |\Sigma| = |AA^T| = |A||A^T| = |A|^2 \\
&= \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} e^{-\frac{1}{2} (x-b)^T \Sigma^{-1} (x-b)}
\end{aligned}$$

□

5.

$$\begin{aligned}
(a) \quad D_{KL}(x || Y) &= \int_{\mathbb{R}^d} f(\alpha) \log\left(\frac{f(\alpha)}{g(\alpha)}\right) d\alpha \\
&= \mathbb{E}\left[-\log \frac{g(\alpha)}{f(\alpha)}\right] \\
&\geq -\log \mathbb{E}\left[\frac{g(\alpha)}{f(\alpha)}\right] \quad (\text{Since } -\log \text{ convex, Jensen's inequality}) \\
&= -\log \int_{\mathbb{R}^d} \frac{g(\alpha)}{f(\alpha)} f(\alpha) d\alpha \\
&= -\log 1 \\
&= 0.
\end{aligned}$$

□

(b)

$$f(\alpha) = \text{pdf}_{X_1}(x_1) \cdots \text{pdf}_{X_d}(x_d)$$

$$g(\underline{\alpha}) = \text{pdf}_{Y_1}(y_1) \cdots \text{pdf}_{Y_d}(y_d)$$

$$\begin{aligned}
D_{KL}(x || Y) &= \int_{\mathbb{R}^d} f(\alpha) \log\left(\frac{f(\alpha)}{g(\alpha)}\right) d\alpha \\
&= \int_{\mathbb{R}} \cdots \int_{\mathbb{R}} \text{pdf}_{X_1}(x_1) \cdots \text{pdf}_{X_d}(x_d) \log \frac{\text{pdf}_{X_1}(x_1) \cdots \text{pdf}_{X_d}(x_d)}{\text{pdf}_{Y_1}(x_1) \cdots \text{pdf}_{Y_d}(x_d)} dx_1 \cdots dx_d \\
&= \sum_{i=1}^d \int_{\mathbb{R}} \cdots \int_{\mathbb{R}} \text{pdf}_{X_1}(x_1) \cdots \text{pdf}_{X_d}(x_d) \log \frac{\text{pdf}_{X_i}(x_i)}{\text{pdf}_{Y_i}(x_i)} dx_1 \cdots dx_d \\
&= \sum_{i=1}^d \int_{\mathbb{R}} \text{pdf}_{X_i}(x_i) \log \frac{\text{pdf}_{X_i}(x_i)}{\text{pdf}_{Y_i}(x_i)} dx_i \\
&= \sum_{i=1}^d D_{KL}(X_i || Y_i)
\end{aligned}$$

□

$$6. \quad X \sim N(\mu_0, \Sigma_0)$$

$$Y \sim N(\mu_1, \Sigma_1)$$

$$\begin{aligned} D_{KL}(X||Y) &= \mathbb{E} \left[-\log \frac{\text{pdf}_Y(x)}{\text{pdf}_X(x)} \right] \\ &= \mathbb{E} \left[-\log \frac{\frac{1}{\sqrt{(2\pi)^d |\Sigma_1|}} \exp \left(-\frac{1}{2} (x-\mu_1)^T \Sigma_1^{-1} (x-\mu_1) \right)}{\frac{1}{\sqrt{(2\pi)^d |\Sigma_0|}} \exp \left(-\frac{1}{2} (x-\mu_0)^T \Sigma_0^{-1} (x-\mu_0) \right)} \right] \\ &= \mathbb{E} \left[-\frac{1}{2} (x-\mu_0)^T \Sigma_0^{-1} (x-\mu_0) \right] + \mathbb{E} \left[\frac{1}{2} (x-\mu_1)^T \Sigma_1^{-1} (x-\mu_1) \right] - \frac{1}{2} \mathbb{E} \left[\log \frac{|\Sigma_0|}{|\Sigma_1|} \right] \end{aligned}$$

Use the fact that $\mathbb{E}(y^T A y) = \text{tr}(A y) + \mu^T A \mu$, where $y \sim N(\mu, \Sigma)$

$$x-\mu_0 \sim N(0, \Sigma_0), \quad x-\mu_1 \sim N(\mu_0-\mu_1, \Sigma_0)$$

$$\begin{aligned} \text{Hence, } D_{KL}(X||Y) &= -\frac{1}{2} \text{tr}(\Sigma_0^{-1} \Sigma_0) + \frac{1}{2} \text{tr}(\Sigma_1^{-1} \Sigma_0) + \frac{1}{2} (\mu_0 - \mu_1)^T \Sigma_1^{-1} (\mu_0 - \mu_1) \\ &\quad + \frac{1}{2} \log \frac{|\Sigma_1|}{|\Sigma_0|} \\ &= \frac{1}{2} \left(\text{tr}(\Sigma_1^{-1} \Sigma_0) + (\mu_1 - \mu_0)^T \Sigma_1^{-1} (\mu_1 - \mu_0) - d + \log \frac{|\Sigma_1|}{|\Sigma_0|} \right) \end{aligned}$$

□