

1 Reversible Architectures [3pts]: In this section, we will investigate a variant for implementing reversible block with affine coupling layers. Consider the following reversible affine coupling block:

$$\begin{aligned} y_1 &= \exp(\mathcal{G}(x_2)) \circ x_1 + \mathcal{F}(x_2) \\ y_2 &= \exp(s) \circ x_2 \end{aligned} \quad (1)$$

where \circ denotes element-wise multiplication. The each inputs $x_1, x_2 \in \mathbb{R}^{\frac{D}{2}}$. The functions \mathcal{F} and \mathcal{G} maps from $\mathbb{R}^{\frac{D}{2}} \rightarrow \mathbb{R}^{\frac{D}{2}}$. This modified block is identical to the ordinary reversible block, except that the inputs x_1 and x_2 are multiplied element-wise by vectors $\exp(\mathcal{F}(x_2))$ and $\exp(s)$.

1. (1pt) Give the equations for inverting this block, i.e. computing x_1 and x_2 from y_2 and y_1 . You may use $/$ to denote element-wise division.

$$\begin{aligned} x_2 &= y_2 / \exp(s) \quad \left(x_{2i} = \frac{y_{2i}}{\exp(s_i)} \right) \\ x_1 &= [y_1 - \mathcal{F}(x_2)] / \exp(\mathcal{G}(x_2)) \\ &= [y_1 - \mathcal{F}(y_2 / \exp(s))] / \exp[\mathcal{G}(y_2 / \exp(s))] \end{aligned}$$

2. (1pt) Give a formula for the Jacobian $\frac{\partial y}{\partial x}$, where y denotes the concatenation of y_1 and y_2 . You may denote the solution as a block matrix, as long as you clearly define what the matrix for each block corresponds to.

$$\frac{\partial y_1}{\partial x_1} = \exp(\mathcal{G}(x_2)), \quad \frac{\partial y_1}{\partial x_2} = \mathcal{F}'(x_2) + \mathcal{G}'(x_2) \circ \exp(\mathcal{G}(x_2)) \circ x_1$$

$$\frac{\partial y_2}{\partial x_1} = 0, \quad \frac{\partial y_2}{\partial x_2} = \exp(s)$$

$$\text{So } J = \begin{bmatrix} \exp(\mathcal{G}(x_2)) & \mathcal{F}'(x_2) + \mathcal{G}'(x_2) \circ \exp(\mathcal{G}(x_2)) \circ x_1 \\ 0 & \exp(s) \end{bmatrix} \in \mathbb{R}^{D \times D}$$

Where all blocks are diagonal $D/2 \times D/2$ matrices with the stated vectors as their diagonals.

Because of the element-wise nature,

$$\text{all } \frac{\partial y_{ki}}{\partial x_{nj}} = 0 \text{ for } i \neq j$$

3. (1pt) Give a formula for the determinant of the Jacobian from previous part, i.e. compute $\det \left(\frac{\partial y}{\partial x} \right)$.
Is this a volume preserving transformation? Justify your answer.

J is upper triangular, so $|J| = \prod_{i=1}^D J_{ii}$

$$\Rightarrow |J| = \prod_{i=1}^{D/2} \exp(g(x_{2i}) + s_i)$$

$$= \exp\left(\sum_{i=1}^{D/2} g(x_{2i}) + s_i\right) \neq 1$$

unless the sum is 0, which we can assume is not the case.

So since $|J| \neq 1$, this transformation is not volume preserving!

2 Variational Free Energy [6pts]: In this question you will derive some expressions related to variational free energy which is maximized to train a VAE. Recall that the VFE is defined as:

$$\mathcal{F}(q) = \mathbb{E}_q[\log p(\mathbf{x}|\mathbf{z})] - D_{KL}(q(\mathbf{z})||p(\mathbf{z}))$$

where KL divergence is defined as

$$D_{KL}(q(\mathbf{z})||p(\mathbf{z})) = \mathbb{E}_q[\log q(\mathbf{z}) - \log p(\mathbf{z})]$$

We will assume that the prior \mathbf{z} is a standard Gaussian:

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I}) = \prod_{i=1}^D p_i(z_i) = \prod_{i=1}^D \mathcal{N}(z_i; 0, 1)$$

Similarly we will assume that the variational approximation $q(\mathbf{z})$ is a fully factorized (i.e., diagonal) Gaussian:

$$q(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{i=1}^D q_i(z_i) = \prod_{i=1}^D \mathcal{N}(z_i; \mu_i, \sigma_i)$$

1. (1pt) Show that:

$$\mathcal{F}(q) = \log p(\mathbf{x}) - D_{KL}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x}))$$

Haha, this one feels like a typical physics homework!

$$\begin{aligned} \mathcal{F}(q) &= \mathbb{E}_q[\log p(\mathbf{x}|\mathbf{z})] - D_{KL}(q(\mathbf{z})||p(\mathbf{z})) \\ &= \mathbb{E}_q[\log p(\mathbf{x}|\mathbf{z}) - \log q(\mathbf{z}) + \log p(\mathbf{z})] \\ &= \mathbb{E}_q[\log(p(\mathbf{x}|\mathbf{z})p(\mathbf{z})) - \log q(\mathbf{z})] \\ &= \mathbb{E}_q[\log(p(\mathbf{z}|\mathbf{x})p(\mathbf{x})) - \log q(\mathbf{z})] \\ &= \mathbb{E}_q[\log p(\mathbf{x}) + \log p(\mathbf{z}|\mathbf{x}) - \log q(\mathbf{z})] \\ &= \log p(\mathbf{x}) - \mathbb{E}_q[\log q(\mathbf{z}) - \log p(\mathbf{z}|\mathbf{x})] \\ &= \log p(\mathbf{x}) - D_{KL}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x})) \quad \square \end{aligned}$$

2. (1pt) Show that the KL term decomposes as a sum of KL terms for individual dimensions. In particular,

$$D_{KL}(q(\mathbf{z})||p(\mathbf{z})) = \sum_i D_{KL}(q_i(z_i)||p_i(z_i))$$

$$\begin{aligned} D_{KL}(q(\mathbf{z})||p(\mathbf{z})) &= \mathbb{E}_q[\log q(\mathbf{z}) - \log p(\mathbf{z})] \\ &= \mathbb{E}_q[\log\left(\prod_{i=1}^D q_i(z_i)\right) - \log\left(\prod_{i=1}^D p_i(z_i)\right)] \\ &= \mathbb{E}_q\left[\sum_{i=1}^D (\log q_i(z_i) - \log p_i(z_i))\right] \\ &= \sum_{i=1}^D \mathbb{E}_q[\log q_i(z_i) - \log p_i(z_i)] \\ &= \sum_{i=1}^D D_{KL}(q_i(z_i)||p_i(z_i)) \quad \square \end{aligned}$$

3. (2pts) Give an explicit formula for the KL divergence $D_{KL}(q_i(z_i) || p_i(z_i))$. This should be a mathematical expression involving μ_i and σ_i .

$$D_{KL}(q_i(z_i) || p_i(z_i)) = E_q[\log(q_i(z_i)) - \log(p_i(z_i))]$$

$$= \int_{\mathbb{R}} dz q_i(z) (\log q_i(z) - \log p_i(z))$$

$$= \int_{\mathbb{R}} dz q_i(z) \log \frac{q_i(z)}{p_i(z)}$$

$$= \int_{\mathbb{R}} dz q_i(z) \log \left[\frac{(\sqrt{2\pi}\sigma_i)^{-1} \exp(-\frac{(z-\mu_i)^2}{2\sigma_i^2})}{\sqrt{2\pi}^{-1} \exp(-z^2/2)} \right]$$

$$= \int_{\mathbb{R}} dz q_i(z) \left[\log(\sigma_i^{-1}) + \log\left(\frac{\exp(-\frac{(z-\mu_i)^2}{2\sigma_i^2})}{\exp(-z^2/2)}\right) \right]$$

$$\int_{\mathbb{R}} dz q_i(z) = 1$$

$$= \log(\sigma_i^{-1}) - \frac{1}{2\sigma_i^2} \int_{\mathbb{R}} dz q_i(z) (z-\mu_i)^2 + \frac{1}{2} \int_{\mathbb{R}} dz q_i(z) z^2$$

$$= \log(\sigma_i^{-1}) - \frac{\text{Var}_q[z_i]}{2\sigma_i^2} + \frac{1}{2} E_q[z_i^2]$$

$$\text{Var}_q[z_i] = \sigma_i^2$$

$$E_q[z_i] = \mu_i$$

$$= \log(\sigma_i^{-1}) - \frac{1}{2} + \frac{1}{2}(\sigma_i^2 + \mu_i^2) \quad \text{Var}(x) = E(x^2) - E(x)^2$$

$$= \frac{1}{2} (\mu_i^2 + \sigma_i^2 - (\log(\sigma_i^2) - 1))$$

4. (2pts) One way to do gradient descent on the KL term is to apply the formula from above. Another approach is to compute stochastic gradients using the reparameterization trick:

$$\nabla_{\theta} D_{KL}(q_i(z_i) || p_i(z_i)) = \mathbb{E}_{\epsilon} [\nabla_{\theta} t_i]$$

, where

$$\theta = \begin{bmatrix} \mu_i \\ \sigma_i \end{bmatrix}$$

and

$$z_i = \mu_i + \sigma_i \epsilon_i$$

$$r_i = \log q_i(z_i)$$

$$s_i = \log p_i(z_i)$$

$$t_i = r_i - s_i$$

(2)

Show how to compute a stochastic estimate of $\nabla_{\theta} D_{KL}(q_i(z_i) || p_i(z_i))$ by doing backpropagation on the above equations. You may find it helpful to draw the computation graph.

Let's start with $\bar{\Theta}_i = \frac{\partial t_i}{\partial \theta_i} = [\bar{\mu}_i, \bar{\sigma}_i]$:

$$\bar{t}_i = 1, \quad \bar{r}_i = 1, \quad \bar{s}_i = -1$$

$$\bar{z}_i = \bar{r}_i \frac{\partial r_i}{\partial z_i} + \bar{s}_i \frac{\partial s_i}{\partial z_i} = \frac{\partial}{\partial z_i} (\log q_i(z_i)) - \frac{\partial}{\partial z_i} (\log p_i(z_i))$$

$$= \frac{\partial}{\partial z_i} [\log q_i(z_i) - \log p_i(z_i)]$$

$$\frac{\partial}{\partial z_i} \log q_i(z) = \frac{\partial}{\partial z_i} \left[\log \left[(\sqrt{2\pi} \sigma_i)^{-1} \exp \left(-\frac{(z_i - \mu_i)^2}{2\sigma_i^2} \right) \right] \right]$$

$$= \frac{\partial}{\partial z_i} \left[\log \left((\sqrt{2\pi} \sigma_i)^{-1} \right) - \frac{(z_i - \mu_i)^2}{2\sigma_i^2} \right]$$

$$= \frac{\partial}{\partial z_i} \left[-\frac{(z_i - \mu_i)^2}{2\sigma_i^2} \right] = \frac{\mu_i - z_i}{\sigma_i^2}$$

$$\frac{\partial}{\partial z_i} \log p_i(z_i) = \frac{\partial}{\partial z_i} \left[\log \left[\sqrt{2\pi}^{-1} \exp(-z_i^2/2) \right] \right]$$

$$= -\frac{\partial}{\partial z_i} \frac{z_i^2}{2} = -z_i$$

$$\Rightarrow \bar{z}_i = \frac{\mu_i - z_i}{\sigma_i^2} + z_i = \bar{\mu}_i$$

$$\bar{\sigma}_i = \bar{z}_i \epsilon_i = \epsilon_i \left(\frac{\mu_i - z_i}{\sigma_i^2} + z_i \right)$$

We know that $z_i = \mu_i + \sigma_i \varepsilon_i$, so $\varepsilon_i = \frac{z_i - \mu_i}{\sigma_i}$
and thus $\varepsilon_i \sim \mathcal{N}(0, 1)$

$$\begin{aligned} \Rightarrow \frac{\partial}{\partial \mu_i} D_{KL}(q_i(z_i) \| p_i(z_i)) &= E_{\varepsilon_i}[\bar{\mu}_i] \\ &= E_{\varepsilon_i} \left[\frac{\mu_i - z_i}{\sigma_i^2} + z_i \right] \\ &= \underbrace{E_{\varepsilon_i}[-\varepsilon_i]}_{=0} + E_{\varepsilon_i}[\mu_i] + \sigma_i \underbrace{E_{\varepsilon_i}[\varepsilon_i]}_{=0} \\ &= \mu_i \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial \sigma_i} D_{KL}(q_i(z_i) \| p_i(z_i)) &= E_{\varepsilon_i}[\bar{\sigma}_i] \\ &= E_{\varepsilon_i} \left[\varepsilon_i \left(\frac{\mu_i - z_i}{\sigma_i^2} + z_i \right) \right] \\ &= E_{\varepsilon_i} \left[-\varepsilon_i^2 / \sigma_i \right] + E_{\varepsilon_i}[\varepsilon_i z_i] \\ &= -\sigma_i^{-1} + E_{\varepsilon_i}[\varepsilon_i \mu_i] + E_{\varepsilon_i}[\sigma_i \varepsilon_i^2] \\ &\quad \quad \quad = \mu_i E_{\varepsilon_i}[\varepsilon_i] = 0 \quad \quad \quad = \sigma_i E_{\varepsilon_i}[\varepsilon_i^2] = \sigma_i \\ &= \sigma_i - \sigma_i^{-1} \end{aligned}$$

$$\begin{aligned} \text{So: } \text{grad}_{\theta_i} D_{KL}(q_i(z_i) \| p_i(z_i)) &= E_{\varepsilon_i}[\text{grad}_{\theta_i} \bar{t}_i] \\ &= \begin{pmatrix} \mu_i \\ \sigma_i - \sigma_i^{-1} \end{pmatrix} \end{aligned}$$

Taking the derivative of the result from 3. gives the exact same gradient!