

Derivation of the Loss in the CDVAE

Antonio Almudévar

$$\log p_{\theta}(\mathbf{x}|\mathbf{c}) = \iint q_{\phi}(\mathbf{z}_l, \mathbf{z}_u|\mathbf{x}) \log p_{\theta}(\mathbf{x}|\mathbf{c}) d\mathbf{z}_l d\mathbf{z}_u \quad (1)$$

$$= \iint q_{\phi}(\mathbf{z}_l, \mathbf{z}_u|\mathbf{x}) \log \frac{p_{\theta}(\mathbf{x}|\mathbf{z}_l, \mathbf{z}_u, \mathbf{c})p_{\theta}(\mathbf{z}_l, \mathbf{z}_u|\mathbf{c})}{p_{\theta}(\mathbf{z}_l, \mathbf{z}_u|\mathbf{c}, \mathbf{x})} d\mathbf{z}_l d\mathbf{z}_u \quad (2)$$

$$= \iint q_{\phi}(\mathbf{z}_l, \mathbf{z}_u|\mathbf{x}) \log \frac{p_{\theta}(\mathbf{x}|\mathbf{z}_l, \mathbf{z}_u, \mathbf{c})p_{\theta}(\mathbf{z}_l, \mathbf{z}_u|\mathbf{c})q_{\phi}(\mathbf{z}_l, \mathbf{z}_u|\mathbf{x})}{p_{\theta}(\mathbf{z}_l, \mathbf{z}_u|\mathbf{c}, \mathbf{x})q_{\phi}(\mathbf{z}_l, \mathbf{z}_u|\mathbf{x})} d\mathbf{z}_l d\mathbf{z}_u \quad (3)$$

$$\begin{aligned} &= \iint q_{\phi}(\mathbf{z}_l, \mathbf{z}_u|\mathbf{x}) \log p_{\theta}(\mathbf{x}|\mathbf{z}_l, \mathbf{z}_u, \mathbf{c}) d\mathbf{z}_l d\mathbf{z}_u \\ &\quad + \iint q_{\phi}(\mathbf{z}_l, \mathbf{z}_u|\mathbf{x}) \log \frac{q_{\phi}(\mathbf{z}_l, \mathbf{z}_u|\mathbf{x})}{p_{\theta}(\mathbf{z}_l, \mathbf{z}_u|\mathbf{c}, \mathbf{x})} d\mathbf{z}_l d\mathbf{z}_u \\ &\quad - \iint q_{\phi}(\mathbf{z}_l|\mathbf{x})q_{\phi}(\mathbf{z}_u|\mathbf{x}) \log \frac{q_{\phi}(\mathbf{z}_l|\mathbf{x})q_{\phi}(\mathbf{z}_u|\mathbf{x})}{p_{\theta}(\mathbf{z}_l|\mathbf{c})p_{\theta}(\mathbf{z}_u)} d\mathbf{z}_l d\mathbf{z}_u \end{aligned} \quad (4)$$

$$\begin{aligned} &= \mathbb{E}_{q_{\phi}(\mathbf{z}_l, \mathbf{z}_u|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z}_l, \mathbf{z}_u, \mathbf{c})] + D_{KL}(q_{\phi}(\mathbf{z}_l, \mathbf{z}_u|\mathbf{x})||p_{\theta}(\mathbf{z}_l, \mathbf{z}_u|\mathbf{c}, \mathbf{x})) \\ &\quad - \int q_{\phi}(\mathbf{z}_u|\mathbf{x}) \int q_{\phi}(\mathbf{z}_l|\mathbf{x}) \log \frac{q_{\phi}(\mathbf{z}_l|\mathbf{x})}{p_{\theta}(\mathbf{z}_l|\mathbf{c})} d\mathbf{z}_l d\mathbf{z}_u \\ &\quad - \int q_{\phi}(\mathbf{z}_l|\mathbf{x}) \int q_{\phi}(\mathbf{z}_u|\mathbf{x}) \log \frac{q_{\phi}(\mathbf{z}_u|\mathbf{x})}{p_{\theta}(\mathbf{z}_u)} d\mathbf{z}_l d\mathbf{z}_u \end{aligned} \quad (5)$$

$$\begin{aligned} &= \mathbb{E}_{q_{\phi}(\mathbf{z}_l, \mathbf{z}_u|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z}_l, \mathbf{z}_u, \mathbf{c})] + D_{KL}(q_{\phi}(\mathbf{z}_l, \mathbf{z}_u|\mathbf{x})||p_{\theta}(\mathbf{z}_l, \mathbf{z}_u|\mathbf{c}, \mathbf{x})) \\ &\quad - D_{KL}(q_{\phi}(\mathbf{z}_l|\mathbf{x})||p_{\theta}(\mathbf{z}_l|\mathbf{c})) - D_{KL}(q_{\phi}(\mathbf{z}_u|\mathbf{x})||p_{\theta}(\mathbf{z}_u)) \end{aligned} \quad (6)$$

Since the term $D_{KL}(q_{\phi}(\mathbf{z}_l, \mathbf{z}_u|\mathbf{x})||p_{\theta}(\mathbf{z}_l, \mathbf{z}_u|\mathbf{c}, \mathbf{x}))$ is positive, we have that an ELBO of $\log p_{\theta}(\mathbf{x}|\mathbf{c})$ is:

$$\begin{aligned} \mathcal{L}_{\theta, \phi}(\mathbf{x}, \mathbf{c}, \mathbf{z}_l, \mathbf{z}_u) &= \\ &\mathbb{E}_{q_{\phi}(\mathbf{z}_l, \mathbf{z}_u|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z}_l, \mathbf{z}_u, \mathbf{c})] - D_{KL}(q_{\phi}(\mathbf{z}_l|\mathbf{x})||p_{\theta}(\mathbf{z}_l|\mathbf{c})) - D_{KL}(q_{\phi}(\mathbf{z}_u|\mathbf{x})||p_{\theta}(\mathbf{z}_u)) \end{aligned} \quad (7)$$

Moreover, since we have imposed that the decoder does not depend on the class, we have that:

$$\mathcal{L}_{CDVAE} = \mathbb{E}_{q_{\phi}(\mathbf{z}_l, \mathbf{z}_u|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z}_l, \mathbf{z}_u)] - D_{KL}(q_{\phi}(\mathbf{z}_l|\mathbf{x})||p_{\theta}(\mathbf{z}_l|\mathbf{c})) - D_{KL}(q_{\phi}(\mathbf{z}_u|\mathbf{x})||p_{\theta}(\mathbf{z}_u)) \quad (8)$$