

# Variational autoencoder and friends

Xun Zheng  
xzheng1@andrew.cmu.edu

## 1 Variational autoencoder and friends

- Generative model (decoder):

$$z \sim P_0(Z) = \mathbf{N}(Z; 0, I) \quad (1)$$

$$x \sim P_\theta(X|z) = \mathbf{N}(X; \mu_\theta(z), I) \quad (2)$$

Generative mutual information:

$$I_{dec}(X; Z) = \mathbf{D}(P_{dec}(X, Z) \| P_{dec}(X)P_{dec}(Z)) \quad (3)$$

where

$$P_{dec}(X, Z) = P_\theta(X|Z)P_0(Z) \quad (4)$$

$$P_{dec}(X) = \int P_{dec}(X, z) dz = \underbrace{\mathbf{E}_{z \sim P_0(Z)} [P_\theta(X|z)]}_{\text{marginal likelihood}} \quad (5)$$

$$P_{dec}(Z) = \int P_{dec}(x, Z) dx = P_0(Z) \quad (6)$$

- Variational posterior takes the form of inference network (encoder):

$$Q_\phi(Z|x) = \mathbf{N}(Z; \mu_\phi(x), \text{diag}(\sigma_\phi^2(x))) \quad (7)$$

Representative mutual information:

$$I_{enc}(X; Z) = \mathbf{D}(P_{enc}(X, Z) \| P_{enc}(X)P_{enc}(Z)) \quad (8)$$

where

$$P_{enc}(X, Z) = Q_\phi(Z|X)P_{data}(X) \quad (9)$$

$$P_{enc}(X) = \int P_{enc}(X, z) dz = P_{data}(X) \quad (10)$$

$$P_{enc}(Z) = \int P_{enc}(x, Z) dx = \underbrace{\mathbf{E}_{x \sim P_{data}(X)} [Q_\phi(Z|x)]}_{\text{aggregated posterior}} \quad (11)$$

- Evidence lower bound for a single sample:

$$\text{ELBO}(x) = \log P_{dec}(x) - \mathbf{D}(Q_\phi(Z|x) \| \underbrace{P_{dec}(Z|x)}_{\text{true posterior}}) \quad (12)$$

$$= \mathbf{E}_{z \sim Q_\phi(Z|x)} \left[ \log \frac{P_{dec}(x, z)}{Q_\phi(z|x)} \right] \quad (13)$$

Minimize the negative ELBO objective:

$$L(\theta, \phi) = \mathbb{E}_{x \sim P_{data}(X)} [-\text{ELBO}(x)] \quad (14)$$

$$= \underbrace{\mathbb{E}_{x \sim P_{data}(X)} \mathbb{E}_{z \sim Q_\phi(Z|x)} [-\log P_\theta(x|z)]}_{\text{negative expected log-likelihood (Nell)}} + \underbrace{\mathbb{E}_{x \sim P_{data}(X)} \mathbb{D}(Q_\phi(Z|x) \| P_0(Z))}_{\text{posterior/prior divergence (Div)}} \quad (15)$$

- In the case of Gaussian prior and posterior, the divergence has a closed form:

$$Div = \mathbb{E}_{x \sim P_{data}(X)} \left[ \frac{1}{2} \cdot \mathbf{1}^T (\sigma_\phi^2(x) + \mu_\phi^2(x) - \mathbf{1} - \log \sigma_\phi^2(x)) \right] \quad (16)$$

- Posterior/prior divergence can be decomposed (Kim and Mnih, 2018):

$$Div = \mathbb{E}_{x \sim P_{data}(X)} \mathbb{E}_{z \sim Q_\phi(Z|x)} \left[ \log \frac{Q_\phi(z|x) P_{enc}(z)}{P_0(z) P_{enc}(z)} \right] \quad (17)$$

$$= \mathbb{E}_{x \sim P_{data}(X)} \mathbb{E}_{z \sim Q_\phi(Z|x)} \left[ \log \frac{Q_\phi(z|x) P_{data}(x)}{P_{enc}(z) P_{data}(x)} \right] + \mathbb{E}_{z \sim P_{enc}(Z)} \left[ \log \frac{P_{enc}(z)}{P_0(z)} \right] \quad (18)$$

$$= I_{enc}(X; Z) + \underbrace{\mathbb{D}(P_{enc}(Z) \| P_0(Z))}_{\text{agg-posterior/prior divergence (Aggdiv)}} \quad (19)$$

- Agg-posterior/prior term above can be further decomposed (Chen et al., 2018; Esmaeili et al., 2019):

$$Aggdiv = \mathbb{E}_{z \sim P_{enc}(Z)} \left[ \log \frac{P_{enc}(z) \prod_j P_{enc}(z_j) \prod_j P_0(z_j)}{P_0(z) \prod_j P_{enc}(z_j) \prod_j P_0(z_j)} \right] \quad (20)$$

$$= \underbrace{\mathbb{D}(P_{enc}(Z) \| \prod_j P_{enc}(Z_j))}_{\text{total correlation (TC)}} + \underbrace{\sum_j \mathbb{D}(P_{enc}(Z_j) \| P_0(Z_j))}_j}_{\text{dimension-wise KL}} - \underbrace{\mathbb{E}_{z \sim P_{enc}(Z)} \left[ \log \frac{P_0(z)}{\prod_j P_0(z_j)} \right]}_{\text{prior factorization}} \quad (21)$$

where  $P_0(Z_j), P_{enc}(Z_j)$  are singleton marginals of  $P_0(Z), P_{enc}(Z)$  respectively. The last term disappears in (Chen et al., 2018) as it assumes independent prior:  $P_0(Z) = \prod_j P_0(Z_j)$ .

- Monte Carlo estimate of agg-posterior entropy:

$$H_{agg}(Z) = \mathbb{E}_{z \sim P_{enc}(Z)} [-\log P_{enc}(z)] \quad (22)$$

$$= \mathbb{E}_{x \sim P_{data}(X)} \mathbb{E}_{z \sim Q_\phi(Z|x)} [-\log \mathbb{E}_{x' \sim P_{data}(X)} Q_\phi(z|x')] \quad (23)$$

- Monte Carlo estimate using full dataset:

$$H_{agg}(Z) \approx \frac{1}{N} \sum_{n=1}^N -\log \left( \frac{1}{N} \sum_{n'=1}^N Q_\phi(z^{(n)} | x^{(n')}) \right), \quad z^{(n)} \sim Q_\phi(Z|x^{(n)}) \quad (24)$$

Cross-evaluate  $N$  posterior samples against  $N$  posterior distributions, hence the complexity is  $\mathcal{O}(N^2)$ . A mini-batch version is straightforward.

- Partially stratified sampling (Esmaeili et al., 2019):

$$H_{agg}(Z) \approx \frac{1}{B} \sum_{b=1}^B -\log P_{enc}(z^{(b)}) \quad (25)$$

Suppose the inner term is estimated using the full dataset, and separate out diagonal ( $n' = b$ ) and off-diagonal terms ( $n' \neq b$ ):

$$P_{enc}(z^{(b)}) \approx \frac{1}{N} \sum_{n'=1}^N Q_\phi(z^{(b)}|x^{(n')}) \quad (26)$$

$$= \frac{1}{N} Q_\phi(z^{(b)}|x^{(b)}) + \frac{1}{N} \sum_{n' \neq b} Q_\phi(z^{(b)}|x^{(n')}) \quad (27)$$

The second term can be estimated from a mini-batch:

$$\frac{1}{N-1} \sum_{n' \neq b} Q_\phi(z^{(b)}|x^{(n')}) \approx \frac{1}{B-1} \sum_{b' \neq b} Q_\phi(z^{(b)}|x^{(b')}) \quad (28)$$

Hence the final estimator is

$$H_{agg}(Z) \approx \frac{1}{|B|} \sum_{b \in B} -\log \left( \frac{1}{N} Q_\phi(z^{(b)}|x^{(b)}) + \frac{N-1}{N} \frac{1}{B-1} \sum_{b' \neq b} Q_\phi(z^{(b)}|x^{(b')}) \right) \quad (29)$$

- Bounds for  $I_{enc}(X; Z)$ :

$$H_{data} - Nell \leq I_{enc}(X; Z) \leq Div \quad (30)$$

Lower bound (Alemi et al., 2018, Appendix D.1):

$$I_{enc}(X; Z) \geq I_{enc}(X; Z) - \mathbb{E}_{z \sim P_{enc}(Z)} \mathbb{D}(P_{enc}(X|z) \| P_\theta(X|z)) \quad (31)$$

$$= \mathbb{E}_{x, z \sim P_{enc}(X, Z)} \left[ \log \frac{P_{enc}(x|z)}{P_{enc}(x)} \right] - \mathbb{E}_{x, z \sim P_{enc}(X, Z)} \left[ \log \frac{P_{enc}(x|z)}{P_\theta(x|z)} \right] \quad (32)$$

$$= H[P_{data}(X)] + \mathbb{E}_{x \sim P_{data}(X)} \mathbb{E}_{z \sim Q_\phi(Z|x)} [\log P_\theta(x|z)] \quad (33)$$

$$= H_{data} - Nell \quad (34)$$

Upper bound (Alemi et al., 2018, Appendix D.2):

$$I_{enc}(X; Z) \leq I_{enc}(X; Z) + Aggdiv = Div \quad (35)$$

- Viewing marginal likelihood as expectation, one can substitute the prior with other distributions using importance sampling trick:

$$P_{dec}(x) = \mathbb{E}_{z \sim P_0(Z)} [P_\theta(x|z)] \quad (36)$$

$$= \mathbb{E}_{z \sim Q_\phi(Z|x)} [P_\theta(x|z)R(x, z)], \quad R(x, z) = \frac{P_0(z)}{Q_\phi(z|x)} \quad (37)$$

$$= \mathbb{E}_{z^{(1)}, \dots, z^{(m)} \sim Q_\phi(Z|x)} \left[ \frac{1}{m} \sum_{i=1}^m P_\theta(x|z^{(i)})R(x, z^{(i)}) \right] \quad (38)$$

The last line may appear to be an unnecessary complication of the second line, however it makes a difference when Jensen's inequality is applied:

$$\log P_{dec}(x) \geq \quad (39)$$

$$\left\{ \underbrace{\mathbb{E}_{z \sim Q_\phi(Z|x)} [\log (P_\theta(x|z)R(x, z))]}_{ELBO(x)} \approx \underbrace{\frac{1}{k} \sum_{j=1}^k \log (P_\theta(x|z^{(j)})R(x, z^{(j)}))}_{MC_k-ELBO(x)} \right. \\ \left. \underbrace{\mathbb{E}_{z^{(1)}, \dots, z^{(m)} \sim Q_\phi(Z|x)} \left[ \log \left( \frac{1}{m} \sum_{i=1}^m P_\theta(x|z^{(i)})R(x, z^{(i)}) \right) \right]}_{IW_m-ELBO(x)} \approx \underbrace{\frac{1}{k} \sum_{j=1}^k \log \left( \frac{1}{m} \sum_{i=1}^m P_\theta(x|z^{(j,i)})R(x, z^{(j,i)}) \right)}_{MC_k-IW_m-ELBO(x)} \right. \quad (40)$$

The quantities are related as follows:

$$\begin{array}{ccccc}
\text{MC}_1\text{-ELBO}(x) & & \text{MC}_1\text{-IW}_m\text{-ELBO}(x) & & \\
\downarrow \approx & \text{often confused when } k=m & \downarrow \approx & & \\
\text{MC}_k\text{-ELBO}(x) & \xrightarrow{=} & \text{MC}_k\text{-IW}_1\text{-ELBO}(x) & & \text{MC}_k\text{-IW}_m\text{-ELBO}(x) \\
\downarrow k \rightarrow \infty & & \downarrow k \rightarrow \infty & & \downarrow k \rightarrow \infty \\
\text{ELBO}(x) & \xrightarrow{=} & \text{IW}_1\text{-ELBO}(x) & \xrightarrow{\leq} & \text{IW}_m\text{-ELBO}(x) \xrightarrow{m \rightarrow \infty} \log P_{dec}(x)
\end{array}$$

Also notice that one can no longer use closed form of  $Div$  in  $\text{IW}_m\text{-ELBO}(x)$  as typically done in VAEs, since log is taken after the averaging over importance samples.

## References

- A. A. Alemi, B. Poole, I. Fischer, J. V. Dillon, R. A. Saurous, and K. Murphy. Fixing a Broken ELBO. In *International Conference on Machine Learning*, 2018.
- R. T. Q. Chen, X. Li, R. Grosse, and D. Duvenaud. Isolating Sources of Disentanglement in Variational Autoencoders. In *Advances in Neural Information Processing Systems*, 2018.
- B. Esmaeili, H. Wu, S. Jain, A. Bozkurt, N. Siddharth, B. Paige, D. H. Brooks, J. Dy, and J.-W. van de Meent. Structured Disentangled Representations. In *International Conference on Artificial Intelligence and Statistics*, 2019.
- H. Kim and A. Mnih. Disentangling by Factorising. In *International Conference on Machine Learning*, 2018.