

IFT6135 – Homework #4

Generative Models

Jean-Philippe Reid
701300

April 14, 2018

1 Reparameterization trick of variational autoencoder

(a)

$$\begin{aligned}\mathbb{E}[z] &= \mathbb{E}[\mu(x) + \sigma(x) \odot \epsilon] \\ &= \mathbb{E}[\mu(x)] + \mathbb{E}[\sigma(x) \odot \epsilon] \\ &= \mu(x) + \mathbb{E}[\sigma(x) \odot \mathcal{N}(\epsilon; 0, \mathbf{I})] \\ &= \mu(x)\end{aligned}\tag{1.1}$$

$$\begin{aligned}\text{Var}[z] &= \text{Var}[\mu(x) + \sigma(x) \odot \epsilon] \\ &= \text{Var}[\mu(x)] + \text{Var}[\sigma(x) \odot \epsilon] \\ &= \text{Var}[\sigma(x) \odot \mathcal{N}(\epsilon; 0, \mathbf{I})] \\ &= \sigma(x)^2\end{aligned}\tag{1.2}$$

Let's suppose that $z = \mu(x) + S(x)\epsilon$, where $S(x) \in \mathbb{R}^{k \times k}$. $\mu(x)$ remains unchanged. The variance is, however, given by

$$\begin{aligned}\text{Var}[S(x)\epsilon] &= S(x) \text{Var}[\epsilon] S^\top(x) \\ &= S(x) \text{Var}[\mathcal{N}(\epsilon; 0, \mathbf{I})] S^\top(x) \\ &= S(x) \mathbf{I} S^\top(x) \\ &= S(x) S^\top(x).\end{aligned}\tag{1.3}$$

Hence, in this specific case, $z \sim \mathcal{N}(\mu(x), S(x) S^\top(x))$, *i.e.* a normal distribution with full covariance matrix given by $S(x) S^\top(x)$.

- (b) Let's first state the main objective of generative models by considering the directed latent-variable model

$$p(x, z | \theta) = p(x | z, \theta) p(z, \theta) = p(z | x, \theta) p(x, \theta). \quad (1.4)$$

The main objective is to learn the parameters θ of p in order to approximate the posterior inference over z given an input x ($p(z | x, \theta)$), but also the marginal inference over x given a latent variable z ($p(x | z, \theta)$). There is, however, a problem with this approach as $p(z | x, \theta)$ is intractable.

The Bayes variational auto-encoder can efficiently solve this problem by maximizing the evidence lower bound (ELBO) defined as

$$\mathcal{L}(p_\theta, q_\phi) = \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x, z) - \log q_\phi(z | x)]. \quad (1.5)$$

The main idea is to find a distribution function $q_\phi(z | x)$ that is the closest to $p_\theta(z | x)$. Then, we will be able to sample directly from the new approximated function $q_\phi(z | x)$ to encode the inputs.

It is important to note that the Eq. 1.5 is directly related to the log likelihood $p_\theta(x)$,

$$\log p_\theta(x) = \mathbb{KL}[q_\phi(z | x) || p_\theta(z | x)] + \mathcal{L}(p_\theta, q_\phi). \quad (1.6)$$

Since x is fixed, we can replace the prior by $q(z | x)$ in Eq. 1.5. This means that $q(z)$ will be different for each input x . Thus, the model will approximate a better posterior $q(z)$.

The optimization problem is twofold:

- (a) optimize over ϕ to keep tight the ELBO around $\log p(x)$
- (b) optimize over θ to increase the lower bound.

Of course, the global effect will be to increase the likelihood $p(x)$.

It is possible to rewrite Eq. 1.6 as

$$\log p_\theta(x) \geq \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \mathbb{KL}[q_\phi(z | x) || p_\theta(z)] \quad (1.7)$$

This equation is the central piece of the variational auto-encoder. Within this framework, $q_\phi(z | x) \sim \mathcal{N}(\mu(x), \text{diag}(\sigma^2(x)))$.

Here are few observations for the variational auto-encoder:

- (a) Both terms on the right hand side of Eq. 1.7 has the term $q_\phi(z | x)$, which can be interpreted as a encoder that describes x .
- (b) In the first term, $\log p_\theta(x | z)$ is the log-likelihood of an input given the latent variable z . In short, this term helps to reconstruct, or decode x given the latent variable z .

- (c) The second term is the \mathbb{KL} divergence between $q_\phi(z | x)$ and the prior $p(z)$, usually fixed to a Gaussian. It prevents $q_\phi(z | x)$ from simply encoding an identity mapping. It forces the model to learn more generalized representations of the data. Hence its appellation as regularization term.
- (d) This task is embedded into a pair of neural networks that are trained to minimize the reconstruction error $|\tilde{x} - x|$ where \tilde{x} is the reconstructed input. The implementation of the reparameterization trick is essential in order to train the model end-to-end.

An alternative approach is to use mean-field inference to resolve this problem. This approach consist on factorizing $q(z | x)$ as

$$q(z | x) = \prod_{j=1}^n q(z_j | x) \quad (1.8)$$

Such representation is based on an important assumption that the latent variables are independent. It helps to approximate all the possible distributions from which we want to inference from.

By putting together the Eqs. 1.5 and 1.8, it is possible to obtain an optimal solution given by

$$q_k(z)^* = \exp \mathbb{E}_{-k} [\log p(z, x)] + \text{cst.} \quad (1.9)$$

The optimal solution is given by the expected value of the joint distribution with respect to all the latent variables except the one we optimized. The usual method to optimize such problem is called *coordinate descent* and consists of optimizing Eqs. 1.5 for the problem in hand for each q_j individually while keeping the other coordinates constant.

Here are few observations for the mean-field inference approach :

- (a) The conditional probability distribution defined in Eq. 1.8 does not contain the *true* posterior since, in reality, the latent variables z_i are dependent with each other. This is also true for the VAE.
- (b) The assumption that $q(z | x)$ is fully factored is too strong considering the rigidity of coordinate descent optimization algorithm. This is in contrast with the variational auto-encoder where all the variables are optimized in a jointly manner through a neural net, $q_\phi(z | x)$. That is the fundamental difference between the two approaches.

In short, the inference variational auto-encoder generally the mean-field method. It can be explained by the capacity of both models to represent the true distribution $q(z | x)$. For the VAE, $q_\phi(z | x)$ which is a function of ϕ . Hence, the the model has

more flexibility to generalize. This is in contrast with the mean field method where $q^{mf}(z \mid x)$ depends on all other latent variable which are set to be constant.

References:

- Lecture notes of CS228, Sections on *The variational auto-encoder* and *Variational inference* – <https://ermongroup.github.io/cs228-notes>
- C.M Bishop, chapter 10.

2 Importance weighted autoencoder

(a) Show that IWLB is a lower bound on the log likelihood $\log p(x)$

$$\begin{aligned} L_k &= \mathbb{E}_{z_{1:k} \sim q(\mathbf{z}|\mathbf{x})} \left[\log \frac{1}{k} \sum_{j=1}^k \frac{p(\mathbf{x}, z_j)}{q(z_j | \mathbf{x})} \right] \\ &\leq \log \left(\mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\frac{1}{k} \sum_{j=1}^k \frac{p(\mathbf{x}, z_j)}{q(z_j | \mathbf{x})} \right] \right) \end{aligned} \quad (2.1)$$

$$\begin{aligned} &= \log \left(\frac{1}{k} \sum_{j=1}^k \mathbb{E}_{q(\mathbf{x}|z_j)} \left[\frac{p(\mathbf{x}, z_j)}{q(z_j | \mathbf{x})} \right] \right) \\ &= \log \left(\frac{1}{k} \sum_{j=1}^k \int q(z_j | \mathbf{x}) \frac{p(\mathbf{x}, z_j)}{q(z_j | \mathbf{x})} d z_j \right) \end{aligned} \quad (2.2)$$

$$\begin{aligned} &= \log \left(\frac{1}{k} \sum_{j=1}^k \int p(\mathbf{x}, z_j) d z_j \right) \\ &= \log \left(\frac{1}{k} \sum_{j=1}^k p(\mathbf{x}) \right) \\ &= \log(p(\mathbf{x})) \end{aligned} \quad (2.3)$$

Eq 2.1 : we used the Jensen's Equality.

Eq 2.2 : The integral can be distributed over each z_j since they are sampled independently.

Hence, more samples can only improve the thightness of the bound!

Reference: Y. Burda *et al.*, Importance weighted autoencoders.

<https://arxiv.org/abs/1509.00519>

(b) L_2 is defined as,

$$L_2 = \mathbb{E}_{z_{1:k} \sim q(\mathbf{z}|\mathbf{x})} \left[\log \frac{1}{2} \sum_{j=1}^2 \frac{p(\mathbf{x}, z_j)}{q(z_j | \mathbf{x})} \right]. \quad (2.4)$$

- $L_2 \leq \log p(x)$

This is a consequence of the Jensen Inequality. See Eq. 2.1 for more details.

•

$$L_2 = \mathbb{E}_{z_{1,2} \sim q(\mathbf{z}|\mathbf{x})} \left[\log \frac{1}{2} \sum_{j=1}^2 \frac{p(\mathbf{x}, z_j)}{q(z_j | \mathbf{x})} \right] \quad (2.5)$$

$$= \mathbb{E}_{z_{1,2} \sim q(\mathbf{z}|\mathbf{x})} \left[\log \mathbb{E}_{z_1 \sim q(\mathbf{z}|\mathbf{x})} \frac{p(\mathbf{x}, z_1)}{q(z_1 | \mathbf{x})} \right] \quad (2.6)$$

$$\geq \mathbb{E}_{z_{1,2} \sim q(\mathbf{z}|\mathbf{x})} \left[\mathbb{E}_{z_1 \sim q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, z_1)}{q(z_1 | \mathbf{x})} \right] \right] \quad (2.7)$$

$$\begin{aligned} &= \mathbb{E}_{z_1 \sim q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, z_1)}{q(z_1 | \mathbf{x})} \right] \\ &= L_1 \end{aligned} \quad (2.8)$$

We used the property of the expectation, $\mathbb{E} \left[1/m \sum_j a_j \right] = 1/n \sum_i a_i$ for the Eq. 2.6; and the Jensen's equality for Eq. 2.7.

Hence, we have that

$$L_1 \leq L_2 \leq \log p(x), \quad (2.9)$$

which shows that L_2 is a tighter bound than ELBO.

3 Maximum likelihood for GANs

We want to find a function $f(\cdot)$ for which the objective function given by

$$J^{(G)} = \max_g \mathbb{E}_{x \sim P_g} [f(D(G(z)))] \quad (3.1)$$

is equivalent to maximizing the likelihood.

It is important to note that maximum likelihood estimation can be interpreted as a process to minimize a KL-Divergence between two distributions. Hence, in order for the cost function $J^{(G)}$ corresponds to the maximum likelihood, we need to show that

$$\frac{\partial}{\partial \theta} J^{(g)} = \frac{\partial}{\partial \theta} \text{KL}(P_d \parallel P_g) = \mathbb{E}_{P_d(x)} \left[\frac{\partial}{\partial \theta} \log P_g(x) \right]. \quad (3.2)$$

where $P_g(x)$ represent the probability distribution of the generated fake image, $G(z)$. Let's develop the left and right hand sides separately,

$$\begin{aligned} \frac{\partial}{\partial \theta} J^{(g)} &= \frac{\partial}{\partial \theta} [\mathbb{E}_{x \sim P_g} [f(D(G(z)))] \\ &= \frac{\partial}{\partial \theta} \left[\int f(D(G(z))) P_g(x) \, dx \right] \\ &= \int f(D(G(z))) \frac{\partial}{\partial \theta} P_g(x) \, dx \\ &= \int f(D(G(z))) P_g(x) \frac{\partial}{\partial \theta} \log P_g(x) \, dx \end{aligned} \quad (3.3)$$

$$\begin{aligned} \text{KL}(P_d(x) \parallel P_g(x)) &= \frac{\partial}{\partial \theta} \mathbb{E}_{P_d(x)} [\log P_d(x) - \log P_g(x)] \\ &= -\frac{\partial}{\partial \theta} \int P_d(x) \log P_g(x) \, dx \\ &= -\int P_d(x) \frac{\partial}{\partial \theta} \log P_g(x) \, dx \end{aligned} \quad (3.4)$$

For Eqs. 3.3 and 3.4 to be equal, we find that

$$f(D(G(z))) = -\frac{P_d(x)}{P_g(x)} \quad (3.5)$$

Then, by assuming that the discriminator is optimal, we find that

$$D^* = \frac{P_d(x)}{P_g(x) + P_d(x)} = \sigma(a(x)) \quad (3.6)$$

$$\begin{aligned}
\frac{P_d(x)}{P_g(x)(1 + P_d(x)/P_g(x))} &= \frac{1}{1 + \exp(a(x))} \\
\frac{-f(x)}{(1 - f(x))} &= \frac{1}{1 + \exp(a(x))} \\
f(x) &= -\exp(a(x))
\end{aligned} \tag{3.7}$$

where we defined $D^* = \sigma(a(x))$ (see Eq. 3.6). This is a valid assumption as $D(x)$ will always be between 0 and 1. By replacing $a(x) = \sigma^{-1}(D^*)$ in Eq. 3.7, we have that

$$f(x) = -\exp(\sigma^{-1}(D^*)). \tag{3.8}$$

For this function, the objective function will be equivalent to maximizing the likelihood.

Reference: Ian Goodfellow, NIPS 2016 Tutorial: Generative Adversarial Networks.
<https://arxiv.org/abs/1701.00160>