

To turn DLVM's intractable learning problem into something tractable, introduce a parametric inference model $g_\phi(z|x)$ which is encoder.

ϕ parameters of inference model; variational parameters.

optimize variational parameters ϕ : $g_\phi(z|x) \approx p_\theta(z|x)$.

Encoder has a directed graphical model, can be factorized:

$$g_\phi(z|x) = g_\phi(z_1, \dots, z_m|x) = \prod_{j=1}^M g_\phi(z_j | \text{Pa}(z_j), x)$$

where $\text{Pa}(z_j)$ is set of parents of z_j in directed graph.

Once we have factorization $g_\phi(z|x)$ can be specifically parameterized using ^(deep) neural networks, where variational parameters ϕ are weights, biases of neural network.

The optimization objective of VAE is evidence lower bound (ELBO), also called variational lower bound.

For any choice of inference model $g_\phi(z|x)$ including choice of ϕ ,

$$\log p_\theta(x) = \mathbb{E}_{g_\phi(z|x)} [\log p_\theta(x)]$$

$$= \mathbb{E}_{g_\phi(z|x)} \left[\log \left(\frac{p_\theta(x, z)}{p_\theta(z|x)} \right) \right]$$

$$= \mathbb{E}_{g_\phi(z|x)} \left[\log \left(\frac{p_\theta(x|z) \cdot g_\phi(z|x)}{g_\phi(z|x) \cdot p_\theta(z|x)} \right) \right]$$

Maximizing ELBO w.r.t. θ, ϕ will concurrently optimize:

- ① approximately maximize marginal likelihood $p_\theta(x)$, so generative model will be better.
- ② minimize KL divergence of approximation $q_\phi(z|x)$ from true posterior $p_\theta(z|x)$ so $q_\phi(z|x)$ becomes better.

Jointly optimize ϕ and θ using SGD.

Reparametrization trick.

for continuous latent variables and differentiable encoder/decoder, the ELBO can be differentiated w.r.t both ϕ and θ if we reparametrize.

First express r.v. $Z \sim q_\phi(z|x)$ as a differentiable/invertible transformation of another r.v. ϵ : Given ϕ, ϵ :

$$Z = g(\epsilon, \phi, x) \quad \text{where } \epsilon \text{ is independent of } x \text{ and } \phi.$$

Now expectations in terms of ϵ :

$$\mathbb{E}_{q_\phi(z|x)} [f(z)] = \mathbb{E}_{p(\epsilon)} [f(z)]$$

and since expectation and gradient commute by linearity:

$$\begin{aligned} \nabla_\phi \mathbb{E}_{q_\phi(z|x)} [f(z)] &= \nabla_\phi \mathbb{E}_{p(\epsilon)} [f(z)] \\ &= \mathbb{E}_{p(\epsilon)} [\nabla_\phi f(z)] \end{aligned}$$

$\simeq \nabla_\phi f(\bar{\epsilon})$ where we can estimate by Monte Carlo draws of $\epsilon \sim p(\epsilon)$ and doing $Z = g(\epsilon, \phi, x)$

(4)

under reparameterization, we can replace expectation w.r.t. $g_\phi(z|x)$ with one w.r.t. $p(\epsilon)$, so ELBO becomes:

$$\begin{aligned} \mathcal{L}_{\theta, \phi}(x) &= \mathbb{E}_{g_\phi(z|x)} [\log p_\theta(x, z) - \log g_\phi(z|x)] \\ &= \mathbb{E}_{p(\epsilon)} [\log p_\theta(x, \epsilon) - \log g_\phi(\epsilon|x)] \end{aligned}$$

$$\text{where } z = g(\epsilon, \phi, x)$$

so we can get simple Monte Carlo estimator $\hat{\mathcal{L}}_{\theta, \phi}(x)$ of individual datapoint ELBO through a single sample $\epsilon \sim p(\epsilon)$

$$\epsilon \sim p(\epsilon)$$

$$z = g(\epsilon, \phi, x)$$

$$\hat{\mathcal{L}}_{\theta, \phi}(x) = \log p_\theta(x, z) - \log g_\phi(z|x)$$

this can be optimized using SGD.

VAEs have some problems:

→ approximate inference distribution is often different from true posterior
(from ELBO objective).

Can modify the ELBO objective itself to balance correct inference
and fitting training data?

InfoVAE

ELBO loss is

$$J_{\text{ELBO}} = -D_{\text{KL}}(q_{\phi}(z) \parallel p(z)) - \mathbb{E}_{q(z)} [D_{\text{KL}}(q_{\phi}(x|z) \parallel p_{\theta}(x|z))]$$

- allow some weighting parameter λ to balance these two terms.
- introduce a mutual information term.

$$J_{\text{InfoVAE}} = -\lambda D_{\text{KL}}(q_{\phi}(z) \parallel p(z)) - \mathbb{E}_{q(z)} [D_{\text{KL}}(q_{\phi}(x|z) \parallel p_{\theta}(x|z))] + \alpha I_g(x; z)$$

where $I_g(x; z)$ is mutual information between x and z under $q_{\phi}(x, z)$

the mutual information maximization encourages the model to use
the latent code and avoids "information preference problem" when
latent variables z ignored in favor of the data variables x .

(6)

$$J_{\text{InfoVAE}} = -\lambda D_{KL}(q_\phi(z) \parallel p(z)) - \mathbb{E}_{q(z)} [D_{KL}(q_\phi(x|z) \parallel p_\theta(x|z))] + \alpha I(x; z)$$

$$= \mathbb{E}_{q_\phi(x, z)} \left[-\lambda \log \frac{q_\phi(z)}{p(z)} - \log \frac{q_\phi(x|z)}{p_\theta(x|z)} - \alpha \log \frac{q_\phi(z)}{q_\phi(z|x)} \right]$$

$$= \mathbb{E}_{q_\phi(x, z)} \left[\log p_\theta(x|z) - \log \frac{q_\phi(z)^{\lambda+\alpha-1} p_\theta(x)}{p(z)^{\lambda} q_\phi(z|x)^{\alpha-1}} \right]$$

$$= \mathbb{E}_{q_\phi(x, z)} \left[\log p_\theta(x|z) - \log \frac{q_\phi(z)^{\lambda+\alpha-1} q_\phi(z|x)^{1-\alpha} p_\theta(x)}{p(z)^{\lambda+\alpha-1} p(z)^{\alpha-1}} \right]$$

$$\boxed{= \mathbb{E}_{p_\theta(x)} \mathbb{E}_{q_\phi(z|x)} \left[\log p_\theta(x|z) \right] - (1-\alpha) \mathbb{E}_{p_\theta(x)} D_{KL}(q_\phi(z|x) \parallel p(z)) - \\ (\alpha+\lambda-1) D_{KL}(q_\phi(z) \parallel p(z)) / -\mathbb{E}_{p_\theta} [\log p_\theta(x)]}$$

↑
constant.