

## **Training Latent Variable Models with Auto-encoding Variational Bayes: A Tutorial**

- Why this tutorial
- Expectation Maximization (EM)
- Approximate E step as variational inference
- Approximate M step
- AEVB
- General recipe for applying AEVB to almost any latent variable models
- Example derivations: Factor Analysis, VAE

Try MLE for latent variable models:

$$\begin{aligned}
 \boldsymbol{\theta}^* &= \arg \max_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(D) \\
 &= \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^N \log p_{\boldsymbol{\theta}}(\mathbf{x}_i) \\
 &= \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^N \log \int p_{\boldsymbol{\theta}}(\mathbf{x}_i, \mathbf{z}_i) d\mathbf{z}_i
 \end{aligned}$$

Generally the integral is intractable: no closed-form solution if  $\mathbf{z}_i$  is continuous; even if  $\mathbf{z}_i$  is discrete there are additional problems that make this less appealing (CAUTION)

One can derive a lower bound to log likelihood using Jensen's inequality:


$$\begin{aligned}
 \sum_{i=1}^N \log p_{\boldsymbol{\theta}}(\mathbf{x}_i) &= \sum_{i=1}^N \log \int q_i(\mathbf{z}_i) \frac{p_{\boldsymbol{\theta}}(\mathbf{x}_i, \mathbf{z}_i)}{q_i(\mathbf{z}_i)} d\mathbf{z}_i \\
 &\geq \sum_{i=1}^N \mathbb{E}_{\mathbf{z}_i \sim q_i(\mathbf{z}_i)} \left[ \log \frac{p_{\boldsymbol{\theta}}(\mathbf{x}_i, \mathbf{z}_i)}{q_i(\mathbf{z}_i)} \right] \quad (\text{equality if } q_i(\mathbf{z}_i) = p_{\boldsymbol{\theta}}(\mathbf{z}_i | \mathbf{x}_i)) \\
 &= \sum_{i=1}^N \mathbb{E}_{\mathbf{z}_i \sim q_i(\mathbf{z}_i)} [\log p_{\boldsymbol{\theta}}(\mathbf{x}_i, \mathbf{z}_i)] + \sum_{i=1}^N \mathbb{H}(q_i)
 \end{aligned}$$

EM algorithm: derive  $q_i(\mathbf{z}_i | \mathbf{x}_i)$  using Bayes rule set  $q_i(\mathbf{z}_i) = \frac{p_{\boldsymbol{\theta}}(\mathbf{z}_i | \mathbf{x}_i)}{p_{\boldsymbol{\theta}}(\mathbf{x}_i)}$  so that bound is tight, then maximize  $\sum_{i=1}^N \mathbb{E}_{\mathbf{z}_i \sim q_i(\mathbf{z}_i)} [\log p_{\boldsymbol{\theta}}(\mathbf{x}_i, \mathbf{z}_i)]$ . At first sight, it might not be clear why maximizing  $\sum_{i=1}^N \mathbb{E}_{\mathbf{z}_i \sim q_i(\mathbf{z}_i)} [\log p_{\boldsymbol{\theta}}(\mathbf{x}_i, \mathbf{z}_i)]$  is easier, but it's indeed easier in practice for plenty of simple models.

$$\begin{aligned}
\sum_{i=1}^N \log p_{\theta}(\mathbf{x}_i) &= \sum_{i=1}^N \mathbb{E}_{\mathbf{z}_i \sim q(\mathbf{z}_i)} [\log p_{\theta}(\mathbf{x}_i)] \\
&= \sum_{i=1}^N \mathbb{E}_{\mathbf{z}_i \sim q(\mathbf{z}_i)} \left[ \log \frac{p_{\theta}(\mathbf{x}_i, \mathbf{z}_i)}{p_{\theta}(\mathbf{z}_i | \mathbf{x}_i)} \right] \\
&= \sum_{i=1}^N \mathbb{E}_{\mathbf{z}_i \sim q(\mathbf{z}_i)} \left[ \log \left( \frac{p_{\theta}(\mathbf{x}_i, \mathbf{z}_i)}{q(\mathbf{z}_i)} \cdot \frac{q(\mathbf{z}_i)}{p_{\theta}(\mathbf{z}_i | \mathbf{x}_i)} \right) \right] \\
&= \underbrace{\sum_{i=1}^N \mathbb{E}_{\mathbf{z}_i \sim q(\mathbf{z}_i)} [\log p_{\theta}(\mathbf{x}_i, \mathbf{z}_i) - \log q(\mathbf{z}_i)]}_{\text{lower bound we had on previously slide}} + D_{\text{KL}}(q(\mathbf{z}_i) \parallel p_{\theta}(\mathbf{z}_i | \mathbf{x}_i))
\end{aligned}$$

We see that the gap between the lower bound and the log likelihood is the sum of all KL-divergences between the chosen distributions  $q(\mathbf{z}_i)$  and the true posteriors!

Goal of E-step: make the lower bound tight

$$\sum_{i=1}^N \log p_{\theta}(\mathbf{x}_i) = \underbrace{\sum_{i=1}^N \mathbb{E}_{\mathbf{z}_i \sim q_i(\mathbf{z}_i)} [\log p_{\theta}(\mathbf{x}_i, \mathbf{z}_i) - \log q_i(\mathbf{z}_i)]}_{\text{lower bound/ELBO}} + \underbrace{\sum_{i=1}^N D_{\text{KL}}(q_i(\mathbf{z}_i) \parallel p_{\theta}(\mathbf{z}_i \mid \mathbf{x}_i))}_{\text{gap}}$$


One way to think about this is to minimize  $D_{\text{KL}}(q_i(\mathbf{z}_i) \parallel p_{\theta}(\mathbf{z}_i \mid \mathbf{x}_i))$  for  $i = 1, \dots, N$ .

But the problem is simply the variational inference problem:

$$q_i^* = \arg \max_{q_i \in \mathcal{Q}} D_{\text{KL}}(q_i(\mathbf{z}_i) \parallel p_{\theta}(\mathbf{z}_i \mid \mathbf{x}_i))$$

How to solve this problem since  $p_{\theta}(\mathbf{z}_i \mid \mathbf{x}_i)$  is **intractable** (or it is tractable but we just want a more general solution that does not require explicit derivations)? There's a solution but two perspectives:

- Since the LHS does not contain  $q_i$ , minimizing the gap with respect to  $q_i$ 's is equivalent to maximizing the ELBO with respect to  $q_i$ . Fortunately, the terms in ELBO are easy to evaluate, and the expectation can be sidestepped with a technique we'll later discuss.
- Apply the "textbook" VI solution:

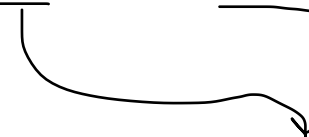
$$D_{\text{KL}}(q_i(\mathbf{z}_i) \parallel \underline{p_{\theta}(\mathbf{z}_i \mid \mathbf{x}_i)}) = D_{\text{KL}}(q_i(\mathbf{z}_i) \parallel \underline{p_{\theta}(\mathbf{x}_i, \mathbf{z}_i)}) + \underbrace{\log p_{\theta}(\mathbf{x}_i)}_{\text{can be dropped}}$$

It is convenient to define  $\underline{Q}$  as a parametrized family of distributions.

But as  $N$  increases, the number of parameters it takes to define all the  $q_i$ 's grow linearly.

It is convenient to represent  $\{q_i\}$  by a neural network  $\underline{q_\phi(\mathbf{z}_i|\mathbf{x}_i)}$  such that  $\underline{q_\phi(\mathbf{z}_i|\mathbf{x}_i)} = q_i(\mathbf{z}_i)$ .

The changes are then reflected in our lower bound:

$$\sum_{i=1}^N \log p_{\theta}(\mathbf{x}_i) \geq \sum_{i=1}^N \mathbb{E}_{\mathbf{z}_i \sim \underline{q_\phi(\mathbf{z}_i|\mathbf{x}_i)}} [\log p_{\theta}(\mathbf{x}_i, \mathbf{z}_i) - \log q_\phi(\mathbf{z}_i|\mathbf{x}_i)]$$


We can maximize ELBO with respect to  $\phi$  with minibatch gradient ascent until convergence:

$$\begin{aligned}
 \phi^{t+1} &\leftarrow \phi^t + \eta \nabla_{\phi} \left\{ \underbrace{\frac{1}{N_B} \sum_{i=1}^{N_B} \mathbb{E}_{\mathbf{z}_i \sim q_{\phi}(\mathbf{z}_i | \mathbf{x}_i)} [\log p_{\theta}(\mathbf{x}_i, \mathbf{z}_i) - \log q_{\phi}(\mathbf{z}_i | \mathbf{x}_i)]}_{\text{ELBO over a minibatch}} \right\}_{\phi = \phi^t} \\
 &= \phi^t + \eta \frac{1}{N_B} \sum_{i=1}^{N_B} \nabla_{\phi} \{ \mathbb{E}_{\mathbf{z}_i \sim q_{\phi}(\mathbf{z}_i | \mathbf{x}_i)} [\log p_{\theta}(\mathbf{x}_i, \mathbf{z}_i) - \log q_{\phi}(\mathbf{z}_i | \mathbf{x}_i)] \}_{\phi = \phi^t} \\
 &= \phi^t + \eta \frac{1}{N_B} \sum_{i=1}^{N_B} \nabla_{\phi} \{ \mathbb{E}_{\epsilon_i \sim q(\epsilon_i)} [\log p_{\theta}(\mathbf{x}_i, \mathbf{z}_i) - \log q_{\phi}(\mathbf{z}_i | \mathbf{x}_i)] \}_{\phi = \phi^t} \quad (\mathbf{z}_i = r(\epsilon_i, \phi, \mathbf{x}_i)) \\
 &= \phi^t + \eta \frac{1}{N_B} \sum_{i=1}^{N_B} \cancel{\mathbb{E}_{\epsilon_i \sim q(\epsilon_i)}} [\nabla_{\phi} \{ \log p_{\theta}(\mathbf{x}_i, \mathbf{z}_i) - \log q_{\phi}(\mathbf{z}_i | \mathbf{x}_i) \}_{\phi = \phi^t}] \\
 &\simeq \phi^t + \eta \frac{1}{N_B} \sum_{i=1}^{N_B} \nabla_{\phi} \{ \log p_{\theta}(\mathbf{x}_i, \mathbf{z}_i) - \log q_{\phi}(\mathbf{z}_i | \mathbf{x}_i) \}_{\phi = \phi^t} \\
 &= \phi^t + \eta \nabla_{\phi} \left\{ \underbrace{\frac{1}{N_B} \sum_{i=1}^{N_B} \log p_{\theta}(\mathbf{x}_i, \mathbf{z}_i) - \log q_{\phi}(\mathbf{z}_i | \mathbf{x}_i)}_{\text{just evaluate this expression in pytorch!}} \right\}_{\phi = \phi^t}
 \end{aligned}$$

where we have applied the reparametrization trick in line 3. It makes the assumption that sampling from  $q_{\phi}(\mathbf{z}_i | \mathbf{x}_i)$  can be made equivalent to (i) first sampling from some base distribution  $\epsilon_i \sim q(\epsilon_i)$  that's free of  $\phi$  and then (ii) transforming the same with a deterministic and differentiable function  $r(\epsilon_i, \phi, \mathbf{x}_i)$ .

Now that ELBO is a tight bound to log likelihood, we can maximize ELBO with respect to generative parameters  $\theta$  with minibatch gradient ascent until convergence:

$$\begin{aligned}
 \theta^{t+1} &\leftarrow \theta^t + \eta \nabla_{\theta} \left\{ \underbrace{\frac{1}{N_B} \sum_{i=1}^{N_B} \mathbb{E}_{\mathbf{z}_i \sim q_{\phi}(\mathbf{z}_i | \mathbf{x}_i)} [\log p_{\theta}(\mathbf{x}_i, \mathbf{z}_i) - \log q_{\phi}(\mathbf{z}_i | \mathbf{x}_i)]}_{\text{ELBO over a minibatch}} \right\}_{\theta=\theta^t} \\
 &= \theta^t + \eta \frac{1}{N_B} \sum_{i=1}^{N_B} \nabla_{\theta} \{ \mathbb{E}_{\mathbf{z}_i \sim q_{\phi}(\mathbf{z}_i | \mathbf{x}_i)} [\log p_{\theta}(\mathbf{x}_i, \mathbf{z}_i) - \log q_{\phi}(\mathbf{z}_i | \mathbf{x}_i)] \}_{\theta=\theta^t} \\
 &= \theta^t + \eta \frac{1}{N_B} \sum_{i=1}^{N_B} \mathbb{E}_{\mathbf{z}_i \sim q_{\phi}(\mathbf{z}_i | \mathbf{x}_i)} [\nabla_{\theta} \{ \log p_{\theta}(\mathbf{x}_i, \mathbf{z}_i) - \log q_{\phi}(\mathbf{z}_i | \mathbf{x}_i) \}_{\theta=\theta^t}] \quad (\text{no reparametrization}) \\
 &\simeq \theta^t + \eta \frac{1}{N_B} \sum_{i=1}^{N_B} \nabla_{\theta} \{ \log p_{\theta}(\mathbf{x}_i, \mathbf{z}_i) - \log q_{\phi}(\mathbf{z}_i | \mathbf{x}_i) \}_{\theta=\theta^t} \quad \text{where} \quad \mathbf{z}_i \sim q_{\phi}(\mathbf{z}_i | \mathbf{x}_i) \\
 &= \theta^t + \eta \nabla_{\theta} \left\{ \underbrace{\frac{1}{N_B} \sum_{i=1}^{N_B} \log p_{\theta}(\mathbf{x}_i, \mathbf{z}_i) - \log q_{\phi}(\mathbf{z}_i | \mathbf{x}_i)}_{\text{just evaluate this expression in pytorch!}} \right\}_{\theta=\theta^t}
 \end{aligned}$$



1. Define the graphical model; write down the distributions for  $p_{\theta}(\mathbf{z}_i)$  and  $p_{\theta}(\mathbf{x}_i|\mathbf{z}_i)$ .
2. Decide on the distribution for  $q_{\phi}(\mathbf{z}_i|\mathbf{x}_i)$ .
3. Derive how to "sample" from the approximate posterior via the reparametrization trick.
4. For that model, implement  $f$ .
5. Alternate between E-steps and M-steps, or just do them simultaneously until convergence:

$$(\phi^{t+1}, \theta^{t+1}) \leftarrow (\phi^t, \theta^t) + \eta \nabla_{(\phi, \theta)} \left\{ \frac{1}{N_B} \sum_{i=1}^{N_B} f(\mathbf{x}_i, \mathbf{z}_i, \theta, \phi) \right\}_{(\phi, \theta) = (\phi^t, \theta^t)} \quad (8)$$

where we've defined  $f$  to be the shorthand for

$$f(\mathbf{x}_i, \mathbf{z}_i, \theta, \phi) = \log p_{\theta}(\mathbf{x}_i, \mathbf{z}_i) - \log q_{\phi}(\mathbf{z}_i|\mathbf{x}_i)$$

Further variance reduction (the reconstruction-KL interpretation):

$$\begin{aligned} & \mathbb{E}_{\mathbf{z}_i \sim q_{\phi}(\mathbf{z}_i|\mathbf{x}_i)} [\log p_{\theta}(\mathbf{x}_i, \mathbf{z}_i) - \log q_{\phi}(\mathbf{z}_i|\mathbf{x}_i)] \\ = & \mathbb{E}_{\mathbf{z}_i \sim q_{\phi}(\mathbf{z}_i|\mathbf{x}_i)} [\log p_{\theta}(\mathbf{x}_i|\mathbf{z}_i) + \log p_{\theta}(\mathbf{z}_i) - \log q_{\phi}(\mathbf{z}_i|\mathbf{x}_i)] \\ = & \mathbb{E}_{\mathbf{z}_i \sim q_{\phi}(\mathbf{z}_i|\mathbf{x}_i)} [\log p_{\theta}(\mathbf{x}_i|\mathbf{z}_i)] - \nabla_{\phi} \mathbb{E}_{\mathbf{z}_i \sim q_{\phi}(\mathbf{z}_i|\mathbf{x}_i)} [\log q_{\phi}(\mathbf{z}_i|\mathbf{x}_i) - \log p_{\theta}(\mathbf{z}_i)] \\ = & \mathbb{E}_{\mathbf{z}_i \sim q_{\phi}(\mathbf{z}_i|\mathbf{x}_i)} [\log p_{\theta}(\mathbf{x}_i|\mathbf{z}_i)] - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}_i|\mathbf{x}_i) || p_{\theta}(\mathbf{z}_i))}_{\text{reconstruction-KL}} \end{aligned}$$

For certain choices of  $q_{\phi}$  and  $p_{\theta}$ , the KL can be evaluated analytically, no need to sample the second term.

(See PDF for code snippets and plots)

(See PDF for code snippets and plots)