

# Variational Inference

Handout Notes - Generative and Deep Learning (GDL)

Davide Bacciu - University of Pisa

---

**Notation.** Random variables are uppercase (e.g.,  $X, Z$ ) and observed values are lowercase (e.g.,  $x, z$ ). A dataset is  $\mathcal{D} = \{x_1, \dots, x_N\}$  with  $N = |\mathcal{D}|$ . Model parameters are  $\theta$ . We use  $p(\cdot)$  for probability mass functions/densities and  $\ell(\theta) = \log p(\mathcal{D} | \theta)$  for the log-likelihood. Latent variables are denoted  $Z$  (and  $z$  for values). The exact posterior is  $p(z | x, \theta)$ . Variational inference introduces a tractable distribution  $q(z | \phi)$  with variational parameters  $\phi$ .

## 1 Why variational inference?

Latent-variable models explain observed data using unobserved variables. A generic model specifies a joint distribution

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

Learning by maximum likelihood aims to maximize

$$\ell(\theta) = \log p(\mathcal{D} | \theta) = \sum_{i=1}^N \log p(x_i | \theta), \quad p(x | \theta) = \int p(x, z | \theta) dz \quad (\text{or } \sum_z \text{ if } Z \text{ is discrete}).$$

The conceptual difficulty is simple: marginalizing out  $z$  often produces integrals/sums without closed form. The computational difficulty is sharper: even when  $p(x, z | \theta)$  is easy to evaluate, the posterior

$$p(z | x, \theta) = \frac{p(x, z | \theta)}{p(x | \theta)}$$

can be intractable because it requires the same marginal likelihood  $p(x | \theta)$  we cannot compute.

Expectation–Maximization (EM) relies on computing posterior expectations in its E-step. When the posterior is intractable, EM in its exact form cannot be applied. Variational inference offers a systematic workaround: replace *sampling/integration* with *optimization* by introducing a tractable approximation  $q(z | \phi)$ .

## 2 From calculus to variational calculus (intuition)

Ordinary calculus optimizes over numbers:  $\max_x f(x)$ . Variational methods optimize over *functions* or *distributions*:

$$\max_q \mathcal{F}(q),$$

where  $\mathcal{F}$  is a functional. In our setting, the object we optimize over is a distribution  $q(z | \phi)$ , usually restricted to a tractable family (e.g., fully factorized distributions).

## 3 Two tools: KL divergence and Jensen's inequality

Variational inference is powered by (i) a way to measure how close two distributions are and (ii) a way to lower-bound a difficult quantity.

### 3.1 Kullback–Leibler divergence

For two distributions  $q(z)$  and  $p(z)$ , the KL divergence is

$$\text{KL}(q\|p) = \mathbb{E}_q \left[ \log \frac{q(z)}{p(z)} \right].$$

It is always nonnegative and equals zero if and only if  $q(z) = p(z)$  almost everywhere. In latent-variable learning, the key comparison is between the variational approximation and the true posterior:

$$\text{KL}(q(z | \phi) \| p(z | x, \theta)) = \mathbb{E}_{q(z|\phi)} \left[ \log \frac{q(z | \phi)}{p(z | x, \theta)} \right].$$

Note that KL is not symmetric, so  $\text{KL}(q\|p) \neq \text{KL}(p\|q)$  in general.

### 3.2 Jensen’s inequality in probabilistic form

If  $f$  is concave (such as  $f = \log$ ), then

$$f(\mathbb{E}[Y]) \geq \mathbb{E}[f(Y)].$$

This inequality becomes useful when we can express a difficult log-likelihood as the log of an expectation under  $q$ .

## 4 Bounding the log-likelihood: the ELBO

The central idea of variational inference is to derive a tractable lower bound on  $\log p(x | \theta)$ . This lower bound is the *evidence lower bound* (ELBO).

#### Worked example — Deriving the ELBO with Jensen’s inequality

Start from the marginal likelihood for a single observation  $x$ :

$$\log p(x | \theta) = \log \int p(x, z | \theta) dz.$$

Introduce any distribution  $q(z | \phi)$  that is nonzero where  $p(x, z | \theta)$  is nonzero and use the identity

$$\int p(x, z | \theta) dz = \int q(z | \phi) \frac{p(x, z | \theta)}{q(z | \phi)} dz = \mathbb{E}_{q(z|\phi)} \left[ \frac{p(x, z | \theta)}{q(z | \phi)} \right].$$

Therefore,

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

where the inequality is Jensen’s (log is concave). Define the ELBO

$$\mathcal{L}(x, \theta, \phi) = \mathbb{E}_{q(z|\phi)}[\log p(x, z | \theta)] - \mathbb{E}_{q(z|\phi)}[\log q(z | \phi)].$$

The second term is the entropy of  $q$ , i.e.,  $H(q) = -\mathbb{E}_q[\log q]$ , so equivalently

$$\mathcal{L}(x, \theta, \phi) = \mathbb{E}_q[\log p(x, z | \theta)] + H(q).$$

For a dataset  $\mathcal{D} = \{x_1, \dots, x_N\}$  with (typically independent) latent variables  $z_1, \dots, z_N$ , a common choice is a factorized variational distribution  $q(\mathbf{z} | \phi) = \prod_{i=1}^N q(z_i | \phi)$ , yielding

$$\log p(\mathcal{D} | \theta) = \sum_{i=1}^N \log p(x_i | \theta) \geq \sum_{i=1}^N \mathcal{L}(x_i, \theta, \phi).$$

## 5 How tight is the bound? ELBO and KL as an exact decomposition

The ELBO is not only a lower bound; it is part of an equality that explains its meaning.

**Worked example — Exact relationship:**  $\log p(x | \theta) = \mathcal{L}(x, \theta, \phi) + \text{KL}(q||p)$

Start from the KL divergence to the true posterior:

$$\text{KL}(q(z | \phi) || p(z | x, \theta)) = \mathbb{E}_q[\log q(z | \phi) - \log p(z | x, \theta)].$$

Use Bayes' rule:  $\log p(z | x, \theta) = \log p(x, z | \theta) - \log p(x | \theta)$ . Substitute:

$$\text{KL}(q||p) = \mathbb{E}_q[\log q(z | \phi) - \log p(x, z | \theta) + \log p(x | \theta)].$$

Rearrange terms (note that  $\log p(x | \theta)$  does not depend on  $z$ ):

$$\log p(x | \theta) = \mathbb{E}_q[\log p(x, z | \theta)] - \mathbb{E}_q[\log q(z | \phi)] + \text{KL}(q(z | \phi) || p(z | x, \theta)).$$

Recognizing the ELBO definition yields the identity

$$\boxed{\log p(x | \theta) = \mathcal{L}(x, \theta, \phi) + \text{KL}(q(z | \phi) || p(z | x, \theta)).}$$

Since  $\text{KL}(\cdot||\cdot) \geq 0$ , the ELBO is indeed a lower bound. The bound is tight iff  $q(z | \phi) = p(z | x, \theta)$ .

This decomposition reveals the operational meaning of variational inference:

*Maximizing the ELBO is equivalent to minimizing the KL divergence between  $q(z | \phi)$  and the true posterior, up to the (unknown) constant  $\log p(x | \theta)$ .*

## 6 Variational inference as a double optimization problem

Variational learning typically solves the coupled optimization

$$\max_{\theta, \phi} \sum_{i=1}^N \mathcal{L}(x_i, \theta, \phi).$$

This is a *double optimization*:

- optimize  $\phi$  so that  $q(z | \phi)$  approximates the posterior well,
- optimize  $\theta$  so that the model assigns high probability to the data (as measured by the bound).

Two equivalent perspectives are useful in practice:

1. **Bound maximization:** directly maximize  $\mathcal{L}(x, \theta, \phi)$  over  $\theta$  and  $\phi$ .
2. **KL minimization:** for fixed  $\theta$ , choose  $\phi$  to minimize  $\text{KL}(q(z | \phi) || p(z | x, \theta))$ .

These are the same because of the exact decomposition  $\log p(x | \theta) = \mathcal{L}(x, \theta, \phi) + \text{KL}(q||p)$ .

## 7 Variational EM: a generalized form of EM

Classical EM alternates between computing an exact posterior over  $z$  (E-step) and maximizing an expected complete-data log-likelihood (M-step). Variational EM keeps the alternating structure but replaces the exact posterior with a tractable approximation.

## 7.1 The alternating updates

At iteration  $k$ :

- **Variational E-step:** update  $\phi$  to improve  $q(z | \phi)$  under the current model parameters  $\theta^{(k)}$ :

$$\phi^{(k)} \in \arg \max_{\phi} \sum_{i=1}^N \mathcal{L}(x_i, \theta^{(k)}, \phi) \quad (\text{equivalently, minimize } \text{KL}(q||p) \text{ for fixed } \theta^{(k)}).$$

- **M-step:** update  $\theta$  to improve the bound under the current variational distribution:

$$\theta^{(k+1)} \in \arg \max_{\theta} \sum_{i=1}^N \mathcal{L}(x_i, \theta, \phi^{(k)}).$$

### Worked example — Recovering exact EM as a special case

If the true posterior  $p(z | x, \theta)$  is tractable, we may choose the variational family to include it and set

$$q(z | \phi^{(k)}) = p(z | x, \theta^{(k)}).$$

Then  $\text{KL}(q||p) = 0$  and the ELBO becomes an equality:

$$\mathcal{L}(x, \theta^{(k)}, \phi^{(k)}) = \log p(x | \theta^{(k)}).$$

In this case, the variational E-step coincides with the classical EM E-step, and variational EM reduces to exact EM.

## 8 Pseudocode: variational EM (general template)

### Worked example — Variational EM algorithm (template)

**Input:** dataset  $\mathcal{D} = \{x_i\}_{i=1}^N$ , model  $p(x, z | \theta)$ , variational family  $\{q(z | \phi)\}$ , tolerance  $\varepsilon$ , max iterations  $K_{\max}$ .

**Output:** parameters  $\theta$  and variational parameters  $\phi$ .

**Initialize:** choose  $\theta^{(0)}$  and  $\phi^{(0)}$ . Set  $k \leftarrow 0$ .

**Repeat until convergence or  $k = K_{\max} - 1$ :**

1. **Variational E-step:** update the variational parameters

$$\phi^{(k)} \leftarrow \arg \max_{\phi} \sum_{i=1}^N \mathcal{L}(x_i, \theta^{(k)}, \phi),$$

using (i) closed-form coordinate updates (if available) or (ii) gradient ascent on the ELBO.

2. **M-step:** update the model parameters

$$\theta^{(k+1)} \leftarrow \arg \max_{\theta} \sum_{i=1}^N \mathcal{L}(x_i, \theta, \phi^{(k)}),$$

often resembling maximum likelihood with latent variables where expectations are taken under  $q(z | \phi^{(k)})$ .

3. **Check convergence:** compute the objective

$$\mathcal{J}^{(k+1)} = \sum_{i=1}^N \mathcal{L}(x_i, \theta^{(k+1)}, \phi^{(k)}),$$

and stop if  $|\mathcal{J}^{(k+1)} - \mathcal{J}^{(k)}| \leq \varepsilon$ . Otherwise set  $k \leftarrow k + 1$ .

**Return**  $\theta^{(k+1)}$  and  $\phi^{(k)}$ .

**Practical note.** In large-scale settings, it is common to use minibatches and stochastic gradient ascent on the ELBO (Stochastic Variational Inference). When optimizing  $\phi$  with gradients, careful reparameterizations or control variates may be needed to reduce gradient variance, depending on the model and variational family.