Variational Inference and Learning I

Fundamentals, mean-field VI, and the EM algorithm

Michael U. Gutmann

Probabilistic Modelling and Reasoning (INFR11134) School of Informatics, The University of Edinburgh

Autumn Semester 2025

Recap

- ► Learning and inference often involves integrals that are hard to compute.
- For example:
 - ► Marginalisation/inference: $p(\mathbf{x}) = \int_{\mathbf{y}} p(\mathbf{x}, \mathbf{y}) d\mathbf{y}$
 - Likelihood in case of unobserved variables: $L(\theta) = p(\mathcal{D}; \theta) = \int_{\mathbf{u}} p(\mathbf{u}, \mathcal{D}; \theta) d\mathbf{u}$
- We here discuss a variational approach to (approximate) inference and learning.

History

Variational methods have a long history, in particular in physics. For example:

- Fermat's principle (1650) to explain the path of light: "light travels between two given points along the path of shortest time" (see e.g. http://www.feynmanlectures.caltech.edu/I_26.html)
- Principle of least action in classical mechanics and beyond (see e.g. http://www.feynmanlectures.caltech.edu/II_19.html)
- Finite elements methods to solve problems in fluid dynamics or civil engineering.

Loosely speaking: the general idea is to frame the original problem in terms of an optimisation problem.

Program

- 1. Preparations
- 2. The variational principle
- 3. Application to inference
- 4. Application to learning

Program

- 1. Preparations
 - Concavity of the logarithm and Jensen's inequality
 - Kullback-Leibler divergence and its properties
- 2. The variational principle
- 3. Application to inference
- 4. Application to learning

log(u) is a concave function

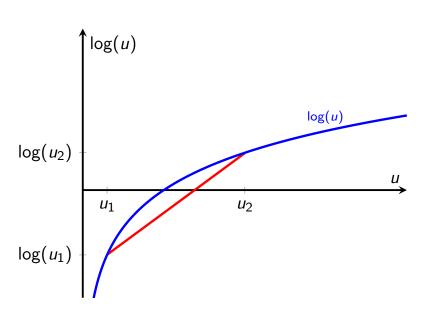
 $ightharpoonup \log(u)$ is a concave function

$$\log((1-a)u_1+au_2)\geq (1-a)\log(u_1)+a\log(u_2)$$
 $a\in [0,1]$ $(1-a)x+ay$ with $a\in [0,1]$ linearly interpolates between x and y .

- ▶ log(average) ≥ average (log)
- Generalisation

$$\log \mathbb{E}[g(\mathbf{x})] \geq \mathbb{E}[\log g(\mathbf{x})]$$

with
$$g(\mathbf{x}) > 0$$



Called Jensen's inequality for concave functions.

Kullback-Leibler divergence

ightharpoonup Kullback Leibler divergence KL(p||q)

$$\mathsf{KL}(p||q) = \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x} = \mathbb{E}_{p(\mathbf{x})} \left[\log \frac{p(\mathbf{x})}{q(\mathbf{x})} \right]$$
(1)

- Properties
 - Arr KL(p||q) = 0 if and only if (iff) p = q (they may be different on sets of probability zero under p)
 - $ightharpoonup \operatorname{\mathsf{KL}}(p||q)
 eq \operatorname{\mathsf{KL}}(q||p)$
 - ightharpoonup KL $(p||q) \ge 0$
- Non-negativity follows from the concavity of the logarithm.

Non-negativity of the KL divergence

Non-negativity follows from the concavity of the logarithm.

$$-\mathsf{KL}(p||q) = -\mathbb{E}_{p(\mathbf{x})} \left[\log \frac{p(\mathbf{x})}{q(\mathbf{x})} \right] \tag{2}$$

$$= \mathbb{E}_{p(\mathbf{x})} \left[\log \frac{q(\mathbf{x})}{p(\mathbf{x})} \right] \tag{3}$$

$$\leq \log \mathbb{E}_{p(\mathbf{x})} \left[\frac{q(\mathbf{x})}{p(\mathbf{x})} \right]$$

$$\int p(\mathbf{x}) \frac{q(\mathbf{x})}{p(\mathbf{x})} d\mathbf{x} = 1$$
(4)

Hence $-\mathsf{KL}(p||q) \leq \log(1) = 0$ and thus

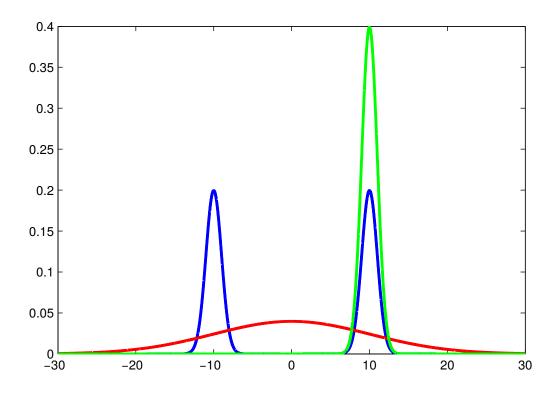
$$\mathsf{KL}(p||q) \ge 0 \tag{5}$$

Asymmetry of the KL divergence

Blue: mixture of Gaussians p(x) (fixed)

Green: (unimodal) Gaussian q that minimises KL(q||p)

Red: (unimodal) Gaussian q that minimises KL(p||q)



Barber Figure 28.1, Section 28.3.4

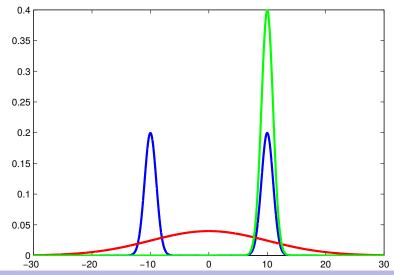
Asymmetry of the KL divergence

$$\operatorname{argmin}_q \mathsf{KL}(q||p) = \operatorname{argmin}_q \int q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})} d\mathbf{x}$$

- ▶ Large penalty when $q(\mathbf{x})$ is large but $p(\mathbf{x})$ is small.
- No penalty when $q(\mathbf{x})$ is small but $p(\mathbf{x})$ is large.
- ▶ Encourages $q(\mathbf{x}) < p(\mathbf{x})$. Produces good local fit, "mode seeking".

$$\operatorname{argmin}_q \mathsf{KL}(p||q) = \operatorname{argmin}_q \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}$$

- ▶ Large penalty when $q(\mathbf{x})$ is small but $p(\mathbf{x})$ is large.
- No penalty when $q(\mathbf{x})$ is large but $p(\mathbf{x})$ is small.
- ▶ Encourages $q(\mathbf{x}) > p(\mathbf{x})$. Produces global fit, corresponds to MLE.

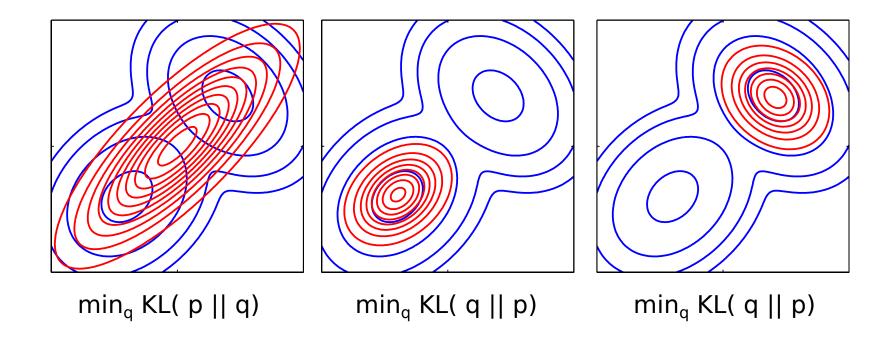


Asymmetry of the KL divergence

Blue: mixture of Gaussians $p(\mathbf{x})$ (fixed)

Red: optimal (unimodal) Gaussians $q(\mathbf{x})$

Global fit (left) versus mode seeking (middle and right). (two local minima are shown)



Bishop Figure 10.3

Program

- 1. Preparations
 - Concavity of the logarithm and Jensen's inequality
 - Kullback-Leibler divergence and its properties
- 2. The variational principle
- 3. Application to inference
- 4. Application to learning

Program

- 1. Preparations
- 2. The variational principle
 - Variational lower bound
 - Maximising the ELBO to compute the marginal and conditional from the joint
- 3. Application to inference
- 4. Application to learning

Variational lower bound: auxiliary distribution

Consider joint pdf /pmf $p(\mathbf{x}, \mathbf{y})$ with marginal $p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{y}) d\mathbf{y}$

ightharpoonup We can write p(x) as

$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{y}) \frac{q(\mathbf{y}|\mathbf{x})}{q(\mathbf{y}|\mathbf{x})} d\mathbf{y} = \mathbb{E}_{q(\mathbf{y}|\mathbf{x})} \left[\frac{p(\mathbf{x}, \mathbf{y})}{q(\mathbf{y}|\mathbf{x})} \right]$$
(6)

where $q(\mathbf{y}|\mathbf{x})$ is an auxiliary distribution (called the variational distribution in the context of variational inference/learning) for a given \mathbf{x} .

Log marginal is

$$\log p(\mathbf{x}) = \log \mathbb{E}_{q(\mathbf{y}|\mathbf{x})} \left[\frac{p(\mathbf{x}, \mathbf{y})}{q(\mathbf{y}|\mathbf{x})} \right]$$
 (7)

➤ Approximating the expectation with a sample average leads to importance sampling. Another approach is to work with the concavity of the logarithm instead.

Variational lower bound: concavity of the logarithm

Concavity of the log gives

$$\log p(\mathbf{x}) = \log \mathbb{E}_{q(\mathbf{y}|\mathbf{x})} \left[\frac{p(\mathbf{x}, \mathbf{y})}{q(\mathbf{y}|\mathbf{x})} \right] \ge \mathbb{E}_{q(\mathbf{y}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{y})}{q(\mathbf{y}|\mathbf{x})} \right]$$
(8)

This is the variational lower bound for $\log p(\mathbf{x})$.

▶ Right-hand side is called the (variational) free energy $\mathcal{F}_{x}(q)$ or the evidence lower bound (ELBO) $\mathcal{L}_{x}(q)$

$$\mathcal{L}_{\mathsf{x}}(q) = \mathbb{E}_{q(\mathsf{y}|\mathsf{x})} \left[\log \frac{p(\mathsf{x}, \mathsf{y})}{q(\mathsf{y}|\mathsf{x})} \right] \tag{9}$$

ightharpoonup Since q is a function, the ELBO is a functional, which is a mapping that depends on a function.

Properties of the ELBO

$$\mathcal{L}_{\mathsf{x}}(q) = \mathbb{E}_{q(\mathsf{y}|\mathsf{x})} \left[\log rac{p(\mathsf{x},\mathsf{y})}{q(\mathsf{y}|\mathsf{x})}
ight]$$

By manipulating the definition of the ELBO, we obtain the following equivalent forms

$$\mathcal{L}_{\mathbf{x}}(q) = \log p(\mathbf{x}) - \mathsf{KL}(q(\mathbf{y}|\mathbf{x})||p(\mathbf{y}|\mathbf{x}))$$
 (10)

$$= \mathbb{E}_{q(\mathbf{y}|\mathbf{x})} \log p(\mathbf{x}|\mathbf{y}) - \mathsf{KL}(q(\mathbf{y}|\mathbf{x})||p(\mathbf{y})) \tag{11}$$

$$= \mathbb{E}_{q(\mathbf{y}|\mathbf{x})} \log p(\mathbf{x}, \mathbf{y}) + \mathcal{H}(q)$$
 (12)

where $p(\mathbf{y})$ is the marginal of $p(\mathbf{x}, \mathbf{y})$ and $\mathcal{H}(q)$ is the entropy of q.

► Entropy is a measure of randomness/variability of a variable

$$\mathcal{H}(q) = -\mathbb{E}_{q(\mathbf{y}|\mathbf{x})} \left[\log q(\mathbf{y}|\mathbf{x}) \right] \tag{13}$$

Larger entropy means more variability.

Properties of the ELBO (proof)

First expression:

$$\mathcal{L}_{\mathbf{x}}(q) = \mathbb{E}_{q(\mathbf{y}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{y})}{q(\mathbf{y}|\mathbf{x})} \right] = \mathbb{E}_{q(\mathbf{y}|\mathbf{x})} \left[\log \frac{p(\mathbf{y}|\mathbf{x})p(\mathbf{x})}{q(\mathbf{y}|\mathbf{x})} \right]$$

$$= \mathbb{E}_{q(\mathbf{y}|\mathbf{x})} \left[\log \frac{p(\mathbf{y}|\mathbf{x})}{q(\mathbf{y}|\mathbf{x})} + \log p(\mathbf{x}) \right]$$

$$= \mathbb{E}_{q(\mathbf{y}|\mathbf{x})} \left[\log \frac{p(\mathbf{y}|\mathbf{x})}{q(\mathbf{y}|\mathbf{x})} + \log p(\mathbf{x}) \right]$$

$$= -KL(q(\mathbf{y}|\mathbf{x})||p(\mathbf{y}|\mathbf{x})) + \log p(\mathbf{x})$$

- Second expression is obtained similarly but using $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x}|\mathbf{y})p(\mathbf{y})$ instead of $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{y}|\mathbf{x})p(\mathbf{x})$ above.
- Third expression from the definition of the entropy.

Tightness of the ELBO

- From $\mathcal{L}_{\mathbf{x}}(q) = \log p(\mathbf{x}) \mathsf{KL}(q(\mathbf{y}|\mathbf{x})||p(\mathbf{y}|\mathbf{x}))$ and non-negativity of the KL divergence, we have
 - 1. $\log p(\mathbf{x}) \geq \mathcal{L}_{\mathbf{x}}(q)$ (as before)
 - 2. $\log p(\mathbf{x}) = \mathcal{L}_{\mathbf{x}}(q) \Leftrightarrow q(\mathbf{y}|\mathbf{x}) = p(\mathbf{y}|\mathbf{x})$
- Maximising $\mathcal{L}_{\mathbf{x}}(q)$ with respect to q yields both $\log p(\mathbf{x})$ and the conditional $p(\mathbf{y}|\mathbf{x})$ at the same time.
- Makes sense: if we know $p(\mathbf{x}, \mathbf{y})$ and $p(\mathbf{x})$, we know $p(\mathbf{y}|\mathbf{x})$, and vice versa, since $p(\mathbf{y}|\mathbf{x}) = p(\mathbf{x}, \mathbf{y})/p(\mathbf{x})$.

Alternative approach

We started from the task of approximating the marginal

$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{y}) d\mathbf{y}$$
 (14)

- Alternative starting point is the task of approximating the conditional $p(\mathbf{y}|\mathbf{x})$ for some given \mathbf{x} by a distribution $q(\mathbf{y}|\mathbf{x})$.
- Measuring the quality of the approximation $q(\mathbf{y}|\mathbf{x})$ by $\mathrm{KL}(q(\mathbf{y}|\mathbf{x})||p(\mathbf{y}|\mathbf{x}))$ gives

$$KL(q(\mathbf{y}|\mathbf{x})||p(\mathbf{y}|\mathbf{x})) = \log p(\mathbf{x}) - \mathcal{L}_{\mathbf{x}}(q)$$
 (15)

Same key result as before.

Variational principle

By maximising the ELBO

$$\mathcal{L}_{\mathbf{x}}(q) = \mathbb{E}_{q(\mathbf{y}|\mathbf{x})} \left[\log rac{p(\mathbf{x}, \mathbf{y})}{q(\mathbf{y}|\mathbf{x})}
ight]$$

we can split the joint $p(\mathbf{x}, \mathbf{y})$ into $p(\mathbf{x})$ and $p(\mathbf{y}|\mathbf{x})$

$$\log p(\mathbf{x}) = \max_{q} \mathcal{L}_{\mathbf{x}}(q)$$
 $p(\mathbf{y}|\mathbf{x}) = \operatorname*{argmax}_{q} \mathcal{L}_{\mathbf{x}}(q)$

► Highlights the variational principle: The inference problem is expressed in terms of an optimisation problem.

Solving the optimisation problem

$$\mathcal{L}_{\mathsf{x}}(q) = \mathbb{E}_{q(\mathsf{y}|\mathsf{x})} \left[\log rac{p(\mathsf{x},\mathsf{y})}{q(\mathsf{y}|\mathsf{x})}
ight]$$

- Difficulties when maximising the ELBO:
 - ▶ Learning of a pdf/pmf $q(\mathbf{y}|\mathbf{x})$
 - Maximisation when objective involves $\mathbb{E}_{q(\mathbf{y}|\mathbf{x})}$ that depends on q
- Restrict search space to a family Q of variational distributions $q(\mathbf{y}|\mathbf{x})$ for which $\mathcal{L}_{\mathbf{x}}(q)$ is computable.
- ightharpoonup Family Q specified by
 - independence assumptions, e.g. $q(\mathbf{y}|\mathbf{x}) = \prod_i q(y_i|\mathbf{x})$, which corresponds to "mean-field" variational inference
 - ▶ parametric assumptions, e.g. $q(y_i|\mathbf{x}) = \mathcal{N}(y_i; \mu_i(\mathbf{x}), \sigma_i^2(\mathbf{x}))$
- Discussed in more detail later.
- \triangleright $\mathcal{L}_{\mathbf{x}}(q)$ can be computed analytically in closed form only in special cases.

Program

- 1. Preparations
- 2. The variational principle
 - Variational lower bound
 - Maximising the ELBO to compute the marginal and conditional from the joint
- 3. Application to inference
- 4. Application to learning

Program

- 1. Preparations
- 2. The variational principle
- 3. Application to inference
 - The mechanics
 - Interpretation
 - Nature of the approximation
 - Mean-field variational inference
- 4. Application to learning

Approximate posterior inference

- Inference task: given value $\mathbf{x} = \mathbf{x}_o$ and joint pdf/pmf $p(\mathbf{x}, \mathbf{y})$, compute $p(\mathbf{y}|\mathbf{x}_o)$.
- Variational approach: estimate the posterior by solving an optimisation problem

$$\hat{p}(\mathbf{y}|\mathbf{x}_o) = \operatorname*{argmax}_{q \in \mathcal{Q}} \mathcal{L}_{\mathbf{x}_o}(q) \tag{16}$$

Q is the set of pdfs/pmfs in which we search for the solution

▶ From the basic property of the ELBO in Equation (10)

$$\log p(\mathbf{x}_o) = \mathsf{KL}(q(\mathbf{y}|\mathbf{x}_o)||p(\mathbf{y}|\mathbf{x}_o)) + \mathcal{L}_{\mathbf{x}_o}(q) = \mathsf{const} \qquad (17)$$

Because the sum of the KL and ELBO is constant, we have

$$\underset{q \in \mathcal{Q}}{\operatorname{argmax}} \, \mathcal{L}_{\mathbf{x}_o}(q) = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} \, \mathsf{KL}(q(\mathbf{y}|\mathbf{x}_o)||p(\mathbf{y}|\mathbf{x}_o)) \tag{18}$$

Posterior as compromise between prior and fit

Equivalent forms of the ELBO:

$$\mathcal{L}_{\mathbf{x}_o}(q) = \mathbb{E}_{q(\mathbf{y}|\mathbf{x}_o)} \log p(\mathbf{x}_o|\mathbf{y}) - \mathsf{KL}(q(\mathbf{y}|\mathbf{x}_o)||p(\mathbf{y})) \tag{19}$$

- ▶ By maximising $\mathcal{L}_{\mathbf{x}_o}(q)$ we find a q that
 - \triangleright produces **y** which are likely explanations of \mathbf{x}_o
 - ightharpoonup stays close to the prior $p(\mathbf{y})$
- ▶ If included in the search space Q, $p(\mathbf{y}|\mathbf{x}_o)$ is the optimal q, which means that the posterior fulfils the two desiderata best.
- ▶ Defines posterior as solution to a regularised decision making problem. (But unlike in the expected loss principle, we here take the expectation with respect to our guess).

As compromise between variable and likely imputations

Equivalent forms of the ELBO:

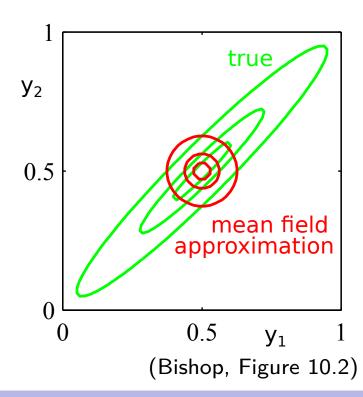
$$\mathcal{L}_{\mathbf{x}_o}(q) = \mathbb{E}_{q(\mathbf{y}|\mathbf{x}_o)} \log p(\mathbf{x}_o, \mathbf{y}) + \mathcal{H}(q)$$
 (20)

- ▶ By maximising $\mathcal{L}_{\mathbf{x}_o}(q)$ we find a q that
 - produces likely imputations (filled-in data) y
 - ► is maximally variable
- ▶ If included in the search space Q, $p(\mathbf{y}|\mathbf{x}_o)$ is the optimal q, which means that the posterior fulfils the two desiderata best.
- ▶ Defines posterior as solution to a regularised decision making problem. (But unlike in the expected loss principle, we here take the expectation with respect to our guess).

Nature of the approximation

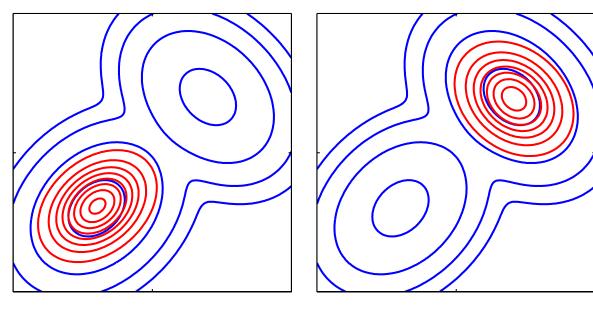
$$\operatorname{argmax}_{q \in \mathcal{Q}} \mathcal{L}_{\mathsf{x}_o}(q) = \operatorname{argmin}_{q \in \mathcal{Q}} \mathsf{KL}(q(\mathsf{y}|\mathsf{x}_o)||p(\mathsf{y}|\mathsf{x}_o))$$

- When minimising KL(q||p) with respect to q, q will try very hard to be zero where p is small.
- Assume true posterior is correlated bivariate Gaussian and we work with $Q = \{q(\mathbf{y}|\mathbf{x}_o) : q(\mathbf{y}|\mathbf{x}_o) = q(y_1|\mathbf{x}_o)q(y_2|\mathbf{x}_o)\}$ (independence but no parametric assumptions)
- ightharpoonup Optimal q is Gaussian.
- Mean is correct but variances dictated by the variances of $p(\mathbf{y}|\mathbf{x}_o)$ along the y_1 and y_2 axes.
- Posterior variance is underestimated.



Nature of the approximation

- Assume that true posterior is multimodal, but that the family of variational distributions \mathcal{Q} only includes unimodal distributions.
- The optimal $q(\mathbf{y}|\mathbf{x}_o)$ only covers one mode: "mode-seeking behaviour".



local optimum

local optimum

Bishop Figure 10.3 (adapted)

Blue: true posterior Red: approximation

Mean-field variational inference

In mean field variational inference, we assume that the variational distribution $q(\mathbf{y}|\mathbf{x})$ fully factorises, i.e.

$$q(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^{d} q_i(y_i|\mathbf{x})$$
 (21)

when \mathbf{y} is d-dimensional.

- ► Independence assumption but no parametric assumption
- An approach to learning the q_i is to update one at a time while keeping the others fixed, called coordinate ascent variational inference (CAVI).
- ▶ We next derive the corresponding update equations.

Mean-field VI: ELBO

With
$$q(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^d q_i(y_i|\mathbf{x})$$
, we have
$$\mathcal{L}_{\mathbf{x}}(q) = \mathbb{E}_{q(\mathbf{y}|\mathbf{x})} \log p(\mathbf{x}, \mathbf{y}) - \mathbb{E}_{q(\mathbf{y}|\mathbf{x})} \log q(\mathbf{y}|\mathbf{x})$$
$$= \mathbb{E}_{q_1(y_1|\mathbf{x})} \cdots \mathbb{E}_{q_d(y_d|\mathbf{x})} \left[\log p(\mathbf{x}, \mathbf{y}) \right] -$$
$$\mathbb{E}_{q_1(y_1|\mathbf{x})} \cdots \mathbb{E}_{q_d(y_d|\mathbf{x})} \left[\sum_{i=1}^d \log q_i(y_i|\mathbf{x}) \right]$$

Second term simplifies:

$$\begin{split} \mathbb{E}_{q_1} \cdots \mathbb{E}_{q_d} \left[\sum_{i=1}^d \log q_i(y_i | \mathbf{x}) \right] &= \sum_{i=1}^d \mathbb{E}_{q_1(y_1 | \mathbf{x})} \cdots \mathbb{E}_{q_d(y_d | \mathbf{x})} \left[\log q_i(y_i | \mathbf{x}) \right] \\ &= \sum_{i=1}^d \mathbb{E}_{q_i(y_i | \mathbf{x})} \left[\log q_i(y_i | \mathbf{x}) \right] \end{split}$$

Mean-field VI: ELBO

Define $q(\mathbf{y}_{\setminus 1}|\mathbf{x}) = \prod_{j=2}^d q_i(y_i|\mathbf{x})$ so that

$$q(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^{d} q_i(y_i|\mathbf{x}) = q_1(y_1|\mathbf{x})q(\mathbf{y}_{\setminus 1}|\mathbf{x})$$
(22)

Hence

$$\mathbb{E}_{q_1(y_1|\mathbf{x})}\cdots\mathbb{E}_{q_d(y_d|\mathbf{x})}\left[\log p(\mathbf{x},\mathbf{y})\right] = \mathbb{E}_{q_1(y_1|\mathbf{x})}\mathbb{E}_{q(\mathbf{y}_{\setminus 1}|\mathbf{x})}\left[\log p(\mathbf{x},\mathbf{y})\right]$$

ELBO becomes

$$\mathcal{L}_{\mathbf{x}}(q) = \mathbb{E}_{q_1(y_1|\mathbf{x})} \mathbb{E}_{q(\mathbf{y}_{\setminus 1}|\mathbf{x})} \left[\log p(\mathbf{x}, \mathbf{y}) \right] - \sum_{i=1}^{d} \mathbb{E}_{q_i(y_i|\mathbf{x})} \left[\log q_i(y_i|\mathbf{x}) \right]$$
(23)

We next maximise $\mathcal{L}_{\mathbf{x}}(q)$ with respect to q_1 while keeping the other q_i fixed.

Mean-field VI: Optimisation

As we optimise with respect to q_1 we can drop additive terms from the ELBO that do not depend on q_1 . This gives the objective

$$J(q_1) = \mathbb{E}_{q_1(y_1|\mathbf{x})} \mathbb{E}_{q(\mathbf{y}_{\setminus 1}|\mathbf{x})} \left[\log p(\mathbf{x}, \mathbf{y}) \right] - \mathbb{E}_{q_1(y_1|\mathbf{x})} \left[\log q_1(y_1|\mathbf{x}) \right]$$
(24)

Define $\bar{p}(y_1|\mathbf{x})$

$$\bar{p}(y_1|\mathbf{x}) = \frac{1}{Z} \exp\left[\mathbb{E}_{q(\mathbf{y}_{\setminus 1}|\mathbf{x})} \left[\log p(\mathbf{x}, \mathbf{y})\right]\right]$$
(25)

Then

$$J(q_1) = \mathbb{E}_{q_1(y_1|\mathbf{x})} \left[\log \bar{p}(y_1|\mathbf{x}) - \log Z \right] - \mathbb{E}_{q_1(y_1|\mathbf{x})} \left[\log q_1(y_1|\mathbf{x}) \right]$$
 (26)

$$= -\mathbb{E}_{q_1(y_1|\mathbf{x})} \left[\log \frac{q_1(y_1|\mathbf{x})}{\bar{p}(y_1|\mathbf{x})} \right] - \log Z \tag{27}$$

$$= -\mathsf{KL}(q_1(y_1|\mathbf{x})||\bar{p}(y_1|\mathbf{x})) - \log Z \tag{28}$$

Thus

$$\underset{q_1}{\operatorname{argmax}} J(q_1) = \underset{q_1}{\operatorname{argmin}} \mathsf{KL}(q_1(y_1|\mathbf{x})||\bar{p}(y_1|\mathbf{x})) = \bar{p}(y_1|\mathbf{x}) \quad (29)$$

Mean-field VI: CAVI updates

► The same calculation for the other marginals gives the CAVI updates:

$$q_i(y_i|\mathbf{x}) = \bar{p}(y_i|\mathbf{x}), \quad \bar{p}(y_i|\mathbf{x}) = \frac{1}{Z} \exp\left[\mathbb{E}_{q(\mathbf{y}_{\setminus i}|\mathbf{x})} \left[\log p(\mathbf{x},\mathbf{y})\right]\right]$$

where $q(\mathbf{y}_{\setminus i}|\mathbf{x}) = \prod_{j \neq i} q(y_j|\mathbf{x})$ is the product of all marginals without marginal $q_i(y_i|\mathbf{x})$.

- Corresponds to coordinate ascent in function space.
- ► In standard CAVI, you update factors in-place and immediately use the newest values for subsequent updates. Guarantees non-decreasing ELBO.
- Order can be predetermined or randomised.

Program

- 1. Preparations
- 2. The variational principle
- 3. Application to inference
 - The mechanics
 - Interpretation
 - Nature of the approximation
 - Mean-field variational inference
- 4. Application to learning

Program

- 1. Preparations
- 2. The variational principle
- 3. Application to inference
- 4. Application to learning
 - Learning with Bayesian models
 - Learning with statistical models and unobserved variables
 - (Variational) EM algorithm

Learning by Bayesian inference

- ▶ Task 1: For a Bayesian model $p(\mathbf{x}|\theta)p(\theta) = p(\mathbf{x},\theta)$, compute the posterior $p(\theta|\mathcal{D})$
- ▶ Formally the same problem as before: $\mathcal{D} = \mathbf{x}_o$ and $\theta \equiv \mathbf{y}$.
- ► Task 2: For a Bayesian model $p(\mathbf{v}, \mathbf{h}|\theta)p(\theta) = p(\mathbf{v}, \mathbf{h}, \theta)$, compute the posterior $p(\theta|\mathcal{D})$ where the data \mathcal{D} are for the visibles \mathbf{v} only.
- ▶ With the equivalence $\mathcal{D} = \mathbf{x}_o$ and $(\mathbf{h}, \boldsymbol{\theta}) \equiv \mathbf{y}$, we are formally back to the problem just studied.

Parameter estimation in presence of unobserved variables

- ► Task: For the model $p(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta})$, estimate the parameters $\boldsymbol{\theta}$ from data \mathcal{D} on the visibles \mathbf{v} only (\mathbf{h} is unobserved).
- ▶ To evaluate the log likelihood function $\ell(\theta)$, we need to evaluate the integral

$$\ell(\boldsymbol{\theta}) = \log p(\mathcal{D}; \boldsymbol{\theta}) = \log \int_{\mathbf{h}} p(\mathcal{D}, \mathbf{h}; \boldsymbol{\theta}) d\mathbf{h}, \tag{30}$$

which is generally intractable.

- In some cases, $\ell(\theta)$ and its gradient can be computed by solving an inference problem, followed by computing an expectation.
- ► Here: use the variational approach.

Parameter estimation in presence of unobserved variables

We had

$$\mathcal{L}_{\mathbf{x}}(q) = \mathbb{E}_{q(\mathbf{y}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{y})}{q(\mathbf{y}|\mathbf{x})} \right]$$

$$= \log p(\mathbf{x}) - \mathsf{KL}(q(\mathbf{y}|\mathbf{x})||p(\mathbf{y}|\mathbf{x}))$$
(31)

Substitute

$$\mathbf{x} \to \mathcal{D}, \quad \mathbf{y} \to \mathbf{h}, \quad p(\mathbf{x}, \mathbf{y}) \to p(\mathcal{D}, \mathbf{h}; \boldsymbol{\theta})$$
 (33)

We then have

$$\mathcal{L}_{\mathcal{D}}(\theta, q) = \mathbb{E}_{q(\mathbf{h}|\mathcal{D})} \left[\log \frac{p(\mathcal{D}, \mathbf{h}; \theta)}{q(\mathbf{h}|\mathcal{D})} \right]$$

$$= \log p(\mathcal{D}; \theta) - \mathsf{KL}(q(\mathbf{h}|\mathcal{D})||p(\mathbf{h}|\mathcal{D}; \theta))$$
(34)

Notation $\mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta},q)$ highlights dependency on $\boldsymbol{\theta}$ and q.

MLE by maximising the ELBO

▶ Using $\ell(\theta)$ for the log-likelihood log $p(\mathcal{D}; \theta)$, we have

$$\mathcal{L}_{\mathcal{D}}(\theta, q) = \ell(\theta) - \mathsf{KL}(q(\mathbf{h}|\mathcal{D})||p(\mathbf{h}|\mathcal{D}; \theta)) \tag{36}$$

▶ If the search space Q is unrestricted or includes $p(\mathbf{h}|\mathcal{D}; \boldsymbol{\theta})$

$$\max_{q} \mathcal{L}_{\mathcal{D}}(\theta, q) = \ell(\theta) \tag{37}$$

Maximum likelihood estimation (MLE)

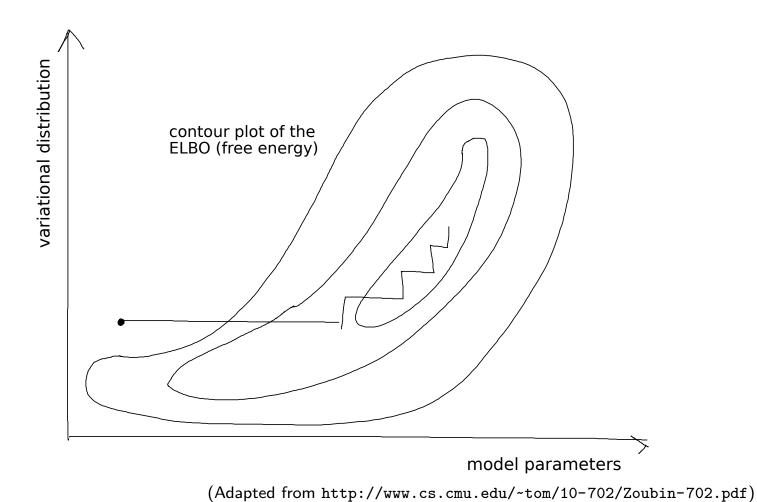
$$\max_{\boldsymbol{\theta}, q} \mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta}, q) = \max_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}) \tag{38}$$

 $\mathsf{MLE} = \mathsf{maximise}$ the ELBO $\mathcal{L}_{\mathcal{D}}(oldsymbol{ heta}, q)$ with respect to $oldsymbol{ heta}$ and q

▶ Restricting the search space Q leads to an approximation when estimating θ and $p(\mathbf{h}|\mathcal{D};\theta)$.

Variational EM algorithm

Variational expectation maximisation (EM): maximise $\mathcal{L}_{\mathcal{D}}(\theta, q)$ by iterating between maximisation with respect to θ and maximisation with respect to q (coordinate ascent).



PMR 2025 ©Gutmann, University of Edinburgh CC BY 4.0

Where is the "expectation"?

► The optimisation with respect to *q* is called the "expectation step"

$$\max_{q \in \mathcal{Q}} \mathcal{L}_{\mathcal{D}}(\theta, q) = \max_{q \in \mathcal{Q}} \mathbb{E}_{q(\mathbf{h}|\mathcal{D})} \left[\log \frac{p(\mathcal{D}, \mathbf{h}; \theta)}{q(\mathbf{h}|\mathcal{D})} \right]$$
(39)

ightharpoonup Denote the best q by q^* so that

$$\max_{q \in \mathcal{Q}} \mathcal{L}_{\mathcal{D}}(\theta, q) = \mathcal{L}_{\mathcal{D}}(\theta, q^*) = \mathbb{E}_{q^*(\mathbf{h}|\mathcal{D})} \left[\log \frac{p(\mathcal{D}, \mathbf{h}; \theta)}{q^*(\mathbf{h}|\mathcal{D})} \right]$$
(40)

which is defined in terms of an expectation and the reason for the name "expectation step".

Classical EM algorithm

- ightharpoonup Denote the parameters at iteration k by θ_k .
- We know that the optimal q for the expectation step is $q^*(\mathbf{h}|\mathcal{D}) = p(\mathbf{h}|\mathcal{D}; \theta_k)$
- If we can compute the posterior $p(\mathbf{h}|\mathcal{D}; \boldsymbol{\theta}_k)$, we obtain the (classical) EM algorithm that iterates between:

E-step: compute the expectation

$$\mathcal{L}_{\mathcal{D}}(\theta, q^*) = \underbrace{\mathbb{E}_{p(\mathbf{h}|\mathcal{D};\theta_k)}[\log p(\mathcal{D}, \mathbf{h}; \theta)]}_{\text{interpretation: expected completed log-likelihood of } - \underbrace{\mathbb{E}_{p(\mathbf{h}|\mathcal{D};\theta_k)}\log p(\mathbf{h}|\mathcal{D}; \theta_k)}_{\text{does not depend on } \theta \text{ and does not need to be computed}}_{\text{does not need to be computed}}$$

M-step: maximise with respect to heta

$$m{ heta}_{k+1} = rgmax_{m{\mathcal{D}}}(m{ heta}, m{q}^*) = rgmax_{m{\mathcal{D}}(m{\mathsf{h}}|m{\mathcal{D}};m{ heta}_k)}[\log p(m{\mathcal{D}}, m{\mathsf{h}};m{ heta})]$$

Classical EM algorithm never decreases the log likelihood

Assume you have updated the parameters and start iteration k+1 with optimisation with respect to q

$$\max_{q} \mathcal{L}_{\mathcal{D}}(\theta_k, q) \tag{41}$$

▶ Optimal solution q_{k+1}^* is the posterior $p(\mathbf{h}|\mathcal{D}; \theta_k)$ so that

$$\ell(\theta_k) = \mathcal{L}_{\mathcal{D}}(\theta_k, q_{k+1}^*) \tag{42}$$

lacktriangle Optimise with respect to the heta while keeping q fixed at q_{k+1}^*

$$\max_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta}, q_{k+1}^*) \tag{43}$$

ightharpoonup Due to maximisation, updated parameter θ_{k+1} is such that

$$\mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta}_{k+1}, q_{k+1}^*) \ge \mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta}_k, q_{k+1}^*) = \ell(\boldsymbol{\theta}_k) \tag{44}$$

From variational lower bound: $\ell(m{ heta}) \geq \mathcal{L}_{\mathcal{D}}(m{ heta}, m{q})$. Hence:

$$\ell(oldsymbol{ heta}_{k+1}) \geq \mathcal{L}_{\mathcal{D}}(oldsymbol{ heta}_{k+1}, q_{k+1}^*) \geq \ell(oldsymbol{ heta}_k)$$

 \Rightarrow EM yields non-decreasing sequence $\ell(\theta_1), \ell(\theta_2), \ldots$

Program recap

1. Preparations

- Concavity of the logarithm and Jensen's inequality
- Kullback-Leibler divergence and its properties

2. The variational principle

- Variational lower bound
- Maximising the ELBO to compute the marginal and conditional from the joint

3. Application to inference

- The mechanics
- Interpretation
- Nature of the approximation
- Mean-field variational inference

4. Application to learning

- Learning with Bayesian models
- Learning with statistical models and unobserved variables
- (Variational) EM algorithm