

Variational Inference

Probabilistic Decision Making — Lecture 9

26th November 2025

Robert Peharz

Institute of Machine Learning and Neural Computation

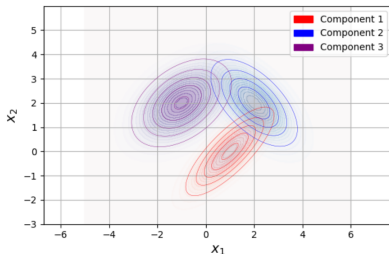
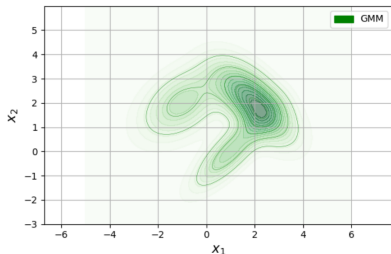
Graz University of Technology

$Z \in \{1, \dots, K\}$ **latent variable**

$$p(\mathbf{x}) = \sum_{k=1}^K w_k p(\mathbf{x} | \mu_k, \Sigma_k)$$

$$p(\mathbf{x}, z) = \underbrace{w_z}_{p(z)} \underbrace{p(\mathbf{x} | \mu_z, \Sigma_z)}_{p(\mathbf{x} | z)}$$

$$p(\mathbf{x}) = \sum_z p(z) p(\mathbf{x} | z)$$



input: Data $\mathcal{D} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$

output: Learned GMM parameters $\theta = \{w_k, \mu_k, \Sigma_k\}_{k=1}^K$

Initialize $\theta = \{w_k, \mu_k, \Sigma_k\}_{k=1}^K$ randomly

for counter = 1 ... max number of iterations **do**

E-Step

for $i = 1 \dots N$, $k = 1 \dots K$ **do**

$$\gamma_{i,k} \leftarrow p(z^{(i)} | \mathbf{x}^{(i)}, \theta) = \frac{p(\mathbf{x}^{(i)} | \mu_k, \Sigma_k) w_k}{\sum_{k'} p(\mathbf{x}^{(i)} | \mu_{k'}, \Sigma_{k'}) w_{k'}} \quad \triangleright \text{responsibilities}$$

end

M-Step

for $k = 1 \dots K$ **do**

$$\mu_k \leftarrow \frac{\sum_{i=1}^N \gamma_{i,k} \mathbf{x}^{(i)}}{\sum_{i=1}^N \gamma_{i,k}} \quad \triangleright \text{weighted mean}$$

$$\Sigma_k \leftarrow \frac{\sum_{i=1}^N \gamma_{i,k} (\mathbf{x}^{(i)} - \mu_k)(\mathbf{x}^{(i)} - \mu_k)^T}{\sum_{i=1}^N \gamma_{i,k}} \quad \triangleright \text{weighted covariance}$$

$$w_k \leftarrow \frac{\sum_{i=1}^N \gamma_{i,k}}{N} \quad \triangleright \text{weighted empirical frequencies}$$

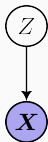
end

end

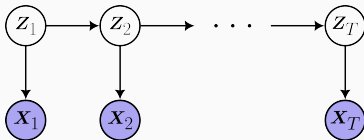
Latent Variable Models

- **latent variables** are a powerful modeling concept in probabilistic ML
- **mixture models** are one example of this concept
- but the concept is much broader, e.g. we can use Bayesian networks with latent variables
- famous example: **hidden Markov models** for sequential data

(Gaussian) Mixture Model



Hidden Markov Model



How to learn models with latent variables or missing data?

EM = maximum likelihood under unobserved data

Given: joint model $p(\mathbf{x}, \mathbf{z} \mid \boldsymbol{\theta})$ over observed \mathbf{X} and unobserved \mathbf{Z}

- **initialize** parameters $\boldsymbol{\theta}_{\text{cur}}$
- **iterate**
 - **step 1** — using the current model, **infer** posterior over latents for each sample $i = 1 \dots N$:

$$p(\mathbf{z}^{(i)} \mid \mathbf{x}^{(i)}, \boldsymbol{\theta}_{\text{cur}})$$

- **step 2** — compute **expected complete log-likelihood**:

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}_{\text{cur}}) = \sum_{i=1}^N \mathbb{E}_{p(\mathbf{z}^{(i)} \mid \mathbf{x}^{(i)}, \boldsymbol{\theta}_{\text{cur}})} \left[\log p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)} \mid \boldsymbol{\theta}) \right]$$

- **step 3** — maximize:

$$\boldsymbol{\theta}_{\text{new}} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}_{\text{cur}}) \text{ and let } \boldsymbol{\theta}_{\text{cur}} \leftarrow \boldsymbol{\theta}_{\text{new}}$$

Why does EM work?

- EM increases monotonically the log-likelihood
- meaning, in each iteration, the marginal log-likelihood can only increase or stay the same—**it can never decrease**
- applicable to missing data in general, not only mixture models
- but, **how and why does EM work?**

Key concept: the **evidence lower bound (ELBO)**, also called **variational lower bound**, **variational objective**, **variational free energy**

Evidence Lower Bound

Setting

- let a parametric joint model $p(\mathbf{x}, \mathbf{z} \mid \theta)$ be given, where
 - \mathbf{X} are observed
 - \mathbf{Z} are unobserved
- **goal:** maximize **marginal log-likelihood** aka **evidence**:

$$\begin{aligned}\arg \max_{\theta} \log p(\mathcal{D} \mid \theta) &= \sum_{i=1}^N \log p(\mathbf{x}^{(i)} \mid \theta) \\ &= \sum_{i=1}^N \log \int p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)} \mid \theta) d\mathbf{z}^{(i)}\end{aligned}$$

- just standard maximum likelihood for latent variable models, i.e. models that are defined as marginal of a bigger model including latent RVs.
- for discrete \mathbf{Z} , the integral is replaced by a sum

In GMMs, each sample has a single latent variable with K states:

$$\begin{aligned}\arg \max_{\boldsymbol{\theta}} \log p(\mathcal{D} | \boldsymbol{\theta}) &= \sum_{i=1}^N \log p(\mathbf{x}^{(i)} | \boldsymbol{\theta}) \\ &= \sum_{i=1}^N \log \sum_{z^{(i)}} p(\mathbf{x}^{(i)}, z^{(i)} | \boldsymbol{\theta}) \\ &= \sum_{i=1}^N \log \sum_{z^{(i)}} w_{z^{(i)}} p(\mathbf{x}^{(i)} | \boldsymbol{\mu}_{z^{(i)}}, \boldsymbol{\Sigma}_{z^{(i)}})\end{aligned}$$

where $\boldsymbol{\theta} = \{w_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K$ are all params of the GMM.

- for each sample i , introduce a **variational distribution** q_i , which can be **any arbitrary** distribution over $\mathbf{Z}^{(i)}$, the unobserved variables in the i^{th} sample
- often, the variational distribution might be **parametric**, but for now we treat them as **non-parametric objects**, without worrying how they are represented
- let $\mathbf{q} = \{q_i\}_{i=1}^N$ be the set of all variational distributions

Evidence Lower Bound (ELBO)

Recall the **marginal log-likelihood** aka **evidence**:

$$\log p(\mathcal{D} | \theta) = \sum_{i=1}^N \overbrace{\log \int p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)} | \theta) d\mathbf{z}^{(i)}}^{\text{sample-wise log-likelihood } \log p(\mathbf{x}^{(i)} | \theta)}$$

We derive the **sample-wise evidence lower bound (ELBO)**:

$$\begin{aligned} \log p(\mathbf{x}^{(i)} | \theta) &= \log \int \overbrace{\frac{q_i(\mathbf{z}^{(i)})}{q_i(\mathbf{z}^{(i)})}}^{=1} p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)} | \theta) d\mathbf{z}^{(i)} \\ &= \log \int q_i(\mathbf{z}^{(i)}) \frac{p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)} | \theta)}{q_i(\mathbf{z}^{(i)})} d\mathbf{z}^{(i)} \\ &= \log \mathbb{E}_{q_i} \left[\frac{p(\mathbf{x}^{(i)}, \mathbf{Z}^{(i)} | \theta)}{q_i(\mathbf{Z}^{(i)})} \right] \\ &\stackrel{\text{Jensen}}{\geq} \mathbb{E}_{q_i} \left[\log \frac{p(\mathbf{x}^{(i)}, \mathbf{Z}^{(i)} | \theta)}{q_i(\mathbf{Z}^{(i)})} \right] =: \text{ELBO}(\theta, q_i) \end{aligned}$$

Interlude: Jensen's Inequality

In the last inequality, we were using

$$\log \mathbb{E}[\dots] \geq \mathbb{E}[\log \dots]$$

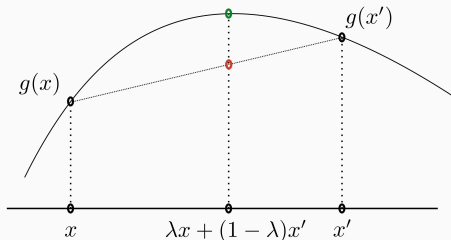
This is due to **Jensen's inequality**: Let \mathbf{X} be any (set of) random variable(s) and g be any **concave function** (e.g. the log). **Then**

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

By the way, if g is a **convex** function then Jensen's inequality reverses:

$$g(\mathbb{E}[\mathbf{X}]) \leq \mathbb{E}[g(\mathbf{X})]$$

$$g(\lambda x + (1 - \lambda)x') \geq \lambda g(x) + (1 - \lambda)g(x')$$



- inequality above is the **definition of concavity**
- $\lambda x + (1 - \lambda)x'$ is a **convex combination** (weighted average) of x and x' , for $\lambda \in [0, 1]$
- **\mathbb{E} operator is also a convex combination!**
- Jensen's inequality: definition of concavity (and convexity) generalizes to expectations

Evidence Lower Bound (ELBO) cont'd

Sample-wise ELBO:

$$\log p(\mathbf{x}^{(i)} | \boldsymbol{\theta}) \geq \mathbb{E}_{q_i} \left[\log \frac{p(\mathbf{x}^{(i)}, \mathbf{Z}^{(i)} | \boldsymbol{\theta})}{q_i(\mathbf{Z}^{(i)})} \right] =: \text{ELBO}(\boldsymbol{\theta}, q_i)$$

The **total ELBO** (or just **ELBO**) is

$$\log p(\mathcal{D} | \boldsymbol{\theta}) = \sum_i \log p(\mathbf{x}^{(i)} | \boldsymbol{\theta}) \geq \sum_i \text{ELBO}(\boldsymbol{\theta}, q_i) =: \text{ELBO}(\boldsymbol{\theta}, \mathbf{q})$$

Note that (total) **ELBO** is a function of both model params $\boldsymbol{\theta}$ and the set of all variational distributions $\mathbf{q} = \{q_i\}_{i=1}^N$. Of course, it also depends on the observed data, but this dependency is usually omitted in notation.

What is the ELBO for?

For **any** set of variational distributions $\mathbf{q} = \{q_i\}_{i=1}^N$ the ELBO is a lower bound of the marginal log-likelihood:

$$\log p(\mathcal{D} | \theta) \geq \text{ELBO}(\theta, \mathbf{q})$$

Since we aim to

$$\max_{\theta} \log p(\mathcal{D} | \theta)$$

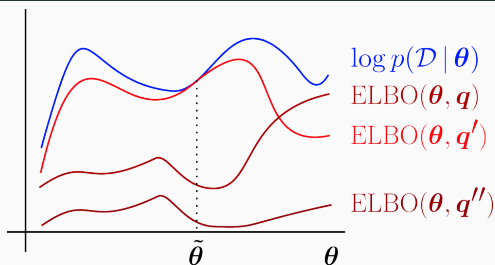
it is also meaningful to

$$\max_{\theta} \text{ELBO}(\theta, \mathbf{q})$$

Maximizing ELBO might sometimes be easier than maximizing the log-likelihood.

Proof for Expectation-Maximization

Illustrating log-likelihood and three ELBOs, for different sets of variational distributions $\mathbf{q}, \mathbf{q}', \mathbf{q}''$



For any \mathbf{q} , it is always true that $\text{ELBO}(\theta, \mathbf{q}) \leq \log p(\mathcal{D} | \theta)$. Moreover, for some \mathbf{q} and some $\tilde{\theta}$, it might happen that

$$\text{ELBO}(\tilde{\theta}, \mathbf{q}) = \log p(\mathcal{D} | \tilde{\theta})$$

In this case we say that the ELBO is **tight** at $\tilde{\theta}$.

In the example above, the ELBO corresponding to \mathbf{q}' is tight at $\tilde{\theta}$.

Interlude: Minorization-Maximization

Forgetting about ELBOs for a moment, **assume** we want to **maximize** some **function**

$$\max_{\mathbf{x}} f(\mathbf{x})$$

but f is “hard to optimize” (not differentiable, expensive to evaluate, highly non-convex, etc.)

Assume that **for each $\tilde{\mathbf{x}}$** we can **construct a function g**

- that **is** a **lower bound**, i.e. for all \mathbf{x} it holds that **$g(\mathbf{x}) \leq f(\mathbf{x})$**
- is **tight** at $\tilde{\mathbf{x}}$, that is **$g(\tilde{\mathbf{x}}) = f(\tilde{\mathbf{x}})$**

Such g is called a **minorization** of f at point $\tilde{\mathbf{x}}$. The idea is to construct a g that is easier to maximize than f .

Interlude: Minorization-Maximization

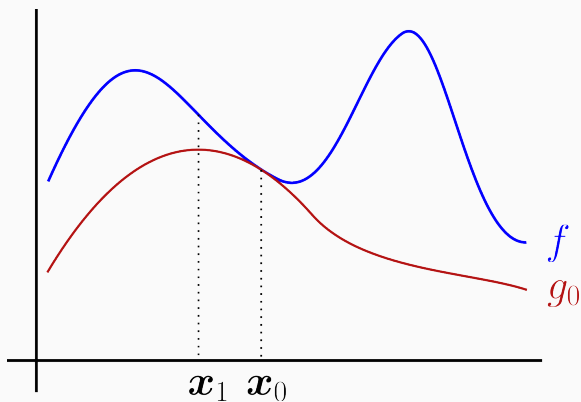
The **minorization-maximization** (MM) principle proceeds as follows:

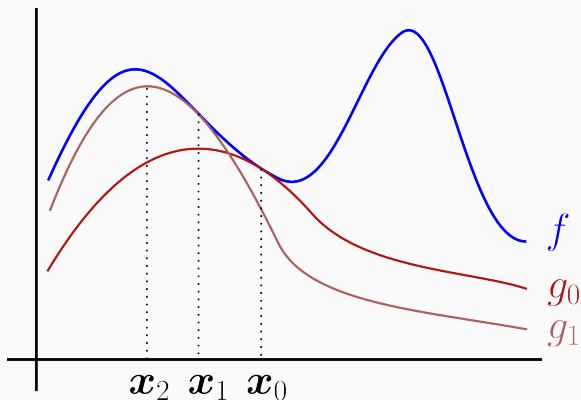
- initialize \mathbf{x}
- iterate:
 - construct minorization g at \mathbf{x}
 - compute $\mathbf{x}^* = \arg \max_{\mathbf{x}} g(\mathbf{x})$
 - $\mathbf{x} \leftarrow \mathbf{x}^*$

It follows

$$f(\mathbf{x}^*) \overset{\text{lower bound}}{\geq} g(\mathbf{x}^*) \overset{\text{maximization}}{\geq} g(\mathbf{x}) \overset{\text{tightness}}{=} f(\mathbf{x})$$

That is, MM monotonically increases f over iterations.

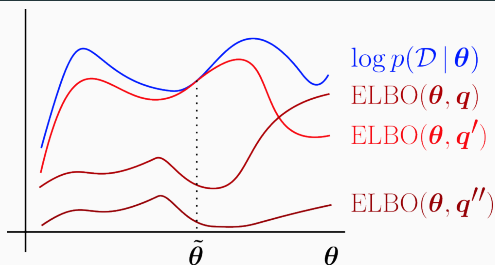




Note: In general, only **local** maximum can be found (MM per se does not even guarantee convergence to local maximum). But, monotonicity over iterations is guaranteed.

ELBO for Minorization-Maximization

Can we use the ELBO for a minorization-maximization scheme?



We'd like to iterate:

- for current model parameters θ_{cur} find q^* such that the ELBO becomes tight at θ_{cur} :

$$\log p(\mathcal{D} | \theta_{\text{cur}}) = \text{ELBO}(\theta_{\text{cur}}, q^*)$$

- $\theta_{\text{new}} = \arg \max_{\theta} \text{ELBO}(\theta, q^*)$
- $\theta_{\text{cur}} \leftarrow \theta_{\text{new}}$

Finding Tight ELBO

Recall the ELBO for the i^{th} sample:

$$\begin{aligned}\text{ELBO}(\theta, q_i) &= \mathbb{E}_{q_i} \left[\log \frac{p(\mathbf{x}^{(i)}, \mathbf{Z}^{(i)} | \theta)}{q_i(\mathbf{Z}^{(i)})} \right] \\ &= \mathbb{E}_{q_i} \left[\log p(\mathbf{x}^{(i)}, \mathbf{Z}^{(i)} | \theta) - \log q_i(\mathbf{Z}^{(i)}) \right]\end{aligned}$$

Using the chain rule, we can write

$$\begin{aligned}p(\mathbf{x}^{(i)}, \mathbf{Z}^{(i)} | \theta) &= p(\mathbf{Z}^{(i)} | \mathbf{x}^{(i)}, \theta) p(\mathbf{x}^{(i)} | \theta) \\ \log p(\mathbf{x}^{(i)}, \mathbf{Z}^{(i)} | \theta) &= \log p(\mathbf{Z}^{(i)} | \mathbf{x}^{(i)}, \theta) + \log p(\mathbf{x}^{(i)} | \theta)\end{aligned}$$

Hence

$$\begin{aligned}\text{ELBO}(\theta, q_i) &= \mathbb{E}_{q_i} \left[\log p(\mathbf{x}^{(i)} | \theta) + \log p(\mathbf{Z}^{(i)} | \mathbf{x}^{(i)}, \theta) - \log q_i(\mathbf{Z}^{(i)}) \right] \\ &= \log p(\mathbf{x}^{(i)} | \theta) + \mathbb{E}_{q_i} \left[\log p(\mathbf{Z}^{(i)} | \mathbf{x}^{(i)}, \theta) - \log q_i(\mathbf{Z}^{(i)}) \right]\end{aligned}$$

Finding Tight ELBO

$$\text{ELBO}(\theta, q_i) = \log p(\mathbf{x}^{(i)} | \theta) + \mathbb{E}_{q_i} \left[\log p(\mathbf{Z}^{(i)} | \mathbf{x}^{(i)}, \theta) - \log q_i(\mathbf{Z}^{(i)}) \right]$$

The second term is

$$\begin{aligned} \mathbb{E}_{q_i} \left[\log p(\mathbf{Z}^{(i)} | \mathbf{x}^{(i)}, \theta) - \log q_i(\mathbf{Z}^{(i)}) \right] &= -\mathbb{E}_{q_i} \left[\log q_i(\mathbf{Z}^{(i)}) - \log p(\mathbf{Z}^{(i)} | \mathbf{x}^{(i)}, \theta) \right] \\ &= -\mathbb{E}_{q_i} \left[\log \frac{q_i(\mathbf{Z}^{(i)})}{p(\mathbf{Z}^{(i)} | \mathbf{x}^{(i)}, \theta)} \right] \\ &= -\text{KL}(q_i || p(\mathbf{z}^{(i)} | \mathbf{x}^{(i)}, \theta)) \end{aligned}$$

that is, the **negative Kullback-Leiber divergence** between the variational distribution q_i and the posterior of $\mathbf{Z}^{(i)}$ given $\mathbf{x}^{(i)}$ under model θ !

Finding Tight ELBO

$$\begin{aligned}\log p(\mathbf{x}^{(i)} | \boldsymbol{\theta}) &= \text{ELBO}(\boldsymbol{\theta}, q_i) + \text{KL}(q_i || p(\mathbf{z}^{(i)} | \mathbf{x}^{(i)}, \boldsymbol{\theta})) \\ \sum_i \log p(\mathbf{x}^{(i)} | \boldsymbol{\theta}) &= \sum_i \text{ELBO}(\boldsymbol{\theta}, q_i) + \sum_i \text{KL}(q_i || p(\mathbf{z}^{(i)} | \mathbf{x}^{(i)}, \boldsymbol{\theta})) \\ \log p(\mathcal{D} | \boldsymbol{\theta}) &= \underbrace{\text{ELBO}(\boldsymbol{\theta}, \mathbf{q}) + \sum_i \text{KL}(q_i || p(\mathbf{z}^{(i)} | \mathbf{x}^{(i)}, \boldsymbol{\theta}))}_{\text{(total) variational gap}}\end{aligned}$$

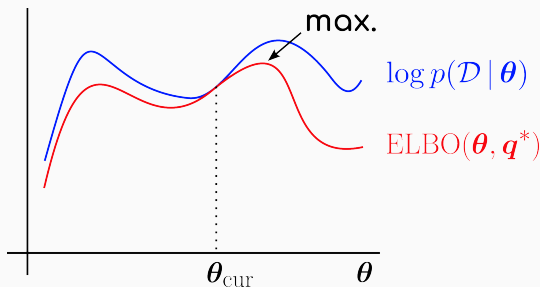
When for each sample q_i is the exact posterior, then $\text{KL} = 0$ and ELBO is tight!

The KL between q_i and $p(\mathbf{z}^{(i)} | \mathbf{x}^{(i)}, \boldsymbol{\theta})$ is called **variational gap** for i^{th} sample. Their sum is the **(total) variational gap**.

$$\log p(\mathcal{D} | \boldsymbol{\theta}) \quad \left| \quad \begin{array}{c} \text{KL} \\ \text{ELBO} \end{array} \right.$$

ELBO for Minorization-Maximization

- let a parameter θ_{cur} of the joint model $p(\mathbf{x}, \mathbf{z} \mid \theta)$ be given
- set $q_i^* \equiv p(\mathbf{z}^{(i)} \mid \mathbf{x}^{(i)}, \theta_{\text{cur}})$ for all i and $\mathbf{q}^* = \{q_i^*\}_{i=1}^N$
- we have established that $\text{ELBO}(\theta, \mathbf{q}^*)$ is a tight lower bound of $\log p(\mathcal{D} \mid \theta)$ at θ_{cur}
- hence, maximizing $\text{ELBO}(\theta, \mathbf{q}^*)$ with respect to θ will increase $\log p(\mathcal{D} \mid \theta)$ (minorization-maximization)



ELBO for Minorization-Maximization cont'd

Goal: $\max_{\theta} \text{ELBO}(\theta, \mathbf{q}^*)$

Recall that

$$\text{ELBO}(\theta, \mathbf{q}^*) = \sum_{i=1}^N \text{ELBO}(\theta, q_i^*)$$

The sample-wise ELBO's can be written as

$$\begin{aligned} \text{ELBO}(\theta, q_i^*) &= \mathbb{E}_{q_i^*} \left[\log \frac{p(\mathbf{x}^{(i)}, \mathbf{Z}^{(i)} | \theta)}{q_i^*(\mathbf{Z}^{(i)})} \right] \\ &= \mathbb{E}_{q_i^*} \left[\log p(\mathbf{x}^{(i)}, \mathbf{Z}^{(i)} | \theta) \right] - \mathbb{E}_{q_i^*} \left[\log q_i^*(\mathbf{Z}^{(i)}) \right] \\ &= \mathbb{E}_{q_i^*} \left[\log p(\mathbf{x}^{(i)}, \mathbf{Z}^{(i)} | \theta) \right] + \underbrace{H[q_i^*]}_{\text{const.}} \end{aligned}$$

Recall that $H[q_i^*]$ is the entropy of q_i^* , which does **not** depend on θ and can be ignored for maximization (it **does** depend on θ_{cur} , which is not optimized though).

ELBO for Minorization-Maximization cont'd

Goal: $\max_{\theta} \text{ELBO}(\theta, \mathbf{q}^*)$

Hence, dropping the entropy, maximizing $\text{ELBO}(\theta, \mathbf{q}^*)$ is equivalent to

$$\max_{\theta} \sum_{i=1}^N \mathbb{E}_{q_i^*} \left[\log p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)} | \theta) \right]$$

But, we have set $q_i^* \equiv p(\mathbf{z}^{(i)} | \mathbf{x}^{(i)}, \theta_{\text{cur}})$ and hence we get

$$\max_{\theta} \sum_{i=1}^N \mathbb{E}_{p(\mathbf{z}^{(i)} | \mathbf{x}^{(i)}, \theta_{\text{cur}})} \left[\log p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)} | \theta) \right]$$

Compare this with slide 5.

ELBO for Minorization-Maximization cont'd

Goal: $\max_{\theta} \text{ELBO}(\theta, \mathbf{q}^*)$

Hence, dropping the entropy, maximizing $\text{ELBO}(\theta, \mathbf{q}^*)$ is equivalent to

$$\max_{\theta} \sum_{i=1}^N \mathbb{E}_{q_i^*} \left[\log p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)} \mid \theta) \right]$$

But, we have set $q_i^* \equiv p(\mathbf{z}^{(i)} \mid \mathbf{x}^{(i)}, \theta_{\text{cur}})$ and hence we get

$$\max_{\theta} \sum_{i=1}^N \mathbb{E}_{p(\mathbf{z}^{(i)} \mid \mathbf{x}^{(i)}, \theta_{\text{cur}})} \left[\log p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)} \mid \theta) \right]$$

Compare this with slide 5. This is exactly EM!

- ELBO is just a lower bound of the log-likelihood, following from Jensen's inequality
- it depends on the variational distributions q_i , hence it is actually a family of lower bounds
- generally, when dealing with maximization problems, maximizing lower bounds is meaningful
- when q_i is the posterior of $\mathbf{Z}^{(i)}$ given $\mathbf{x}^{(i)}$ for a given model θ_{cur} , the ELBO is tight at θ_{cur}
- tight ELBO leads to a minorization-maximization approach, also known as EM

Variational EM

ELBO and Posterior Inference

Recap relation between log-likelihood, ELBO and variational gap:

$$\underbrace{\log p(\mathbf{x}^{(i)} | \boldsymbol{\theta})}_{\text{const. w.r.t. } q_i} = \text{ELBO}(\boldsymbol{\theta}, q_i) + \text{KL}(q_i || p(\mathbf{z}^{(i)} | \mathbf{x}^{(i)}, \boldsymbol{\theta}))$$

- ELBO is a function of both $\boldsymbol{\theta}$ and q_i
- when fixing $\boldsymbol{\theta}$, then $\log p(\mathbf{x}^{(i)} | \boldsymbol{\theta})$ is constant
- hence, increasing $\text{ELBO}(\boldsymbol{\theta}, q_i)$ with respect to q_i will decrease $\text{KL}(q_i || p(\mathbf{z}^{(i)} | \mathbf{x}^{(i)}, \boldsymbol{\theta}))$
- **maximizing ELBO means performing posterior inference!**

ELBO and Posterior Inference cont'd

- global maximum of ELBO at $q_i \equiv p(\mathbf{z}^{(i)}|\mathbf{x}^{(i)}, \boldsymbol{\theta})$
- for GMMs, this is closed form, but for many other models computing the posterior is **hard**
- **Idea:** introduce **parametric variational distributions** $q_{\phi_i}(\mathbf{z}^{(i)})$, where ϕ_i are called **variational parameters**
- **maximizing the ELBO:**

$$\max_{\phi_i} \text{ELBO}(\boldsymbol{\theta}, \phi_i) =: \text{ELBO}(\boldsymbol{\theta}, q_i)$$

is equivalent to minimizing the KL to the posterior:

$$\min_{\phi_i} \text{KL}(q_{\phi_i} || p(\mathbf{z}^{(i)}|\mathbf{x}^{(i)}, \boldsymbol{\theta}))$$

- when the variational distributions do not contain the true posterior, the KL cannot become 0 \Rightarrow **approximate posterior**

Variational EM

- if the variational distribution is restricted, the ELBO might not be tight—nevertheless, the ELBO is always a lower bound of the log-likelihood, so maximizing it is meaningful
- we might optimize the ELBO both with respect to θ and ϕ :

$$\max_{\theta, \phi} \text{ELBO}(\theta, \phi)$$

where $\phi = \{\phi_i\}_{i=1}^N$ is the set of all variational parameters

- maximize w.r.t. ϕ : improve lower bound
- maximize w.r.t. θ : learn joint model with ELBO as a surrogate of the log-likelihood
- this scheme is called **variational EM**
- EM is just a special case where the ELBO is made tight

Let us demonstrate variational EM on GMMs. This is not very meaningful, since exact EM can be done here, but we aim to illustrate the concept.

- variational distributions are naturally Categoricals:

$$q(z^{(i)} = k | \phi_i) = \phi_{i,k}, \quad \text{where } \phi_{i,k} \geq 0, \quad \sum_{k=1}^K \phi_{i,k} = 1$$

- the set of all variational distributions might be interpreted as a non-negative $N \times K$ matrix whose rows sum to one:

$$\phi = \begin{pmatrix} \phi_{1,1} & \phi_{1,2} & \dots & \phi_{1,K} \\ \phi_{2,1} & \phi_{2,2} & \dots & \phi_{2,K} \\ \vdots & \vdots & \ddots & \vdots \\ \phi_{N,1} & \phi_{N,2} & \dots & \phi_{N,K} \end{pmatrix}$$

We reparametrize ϕ by an unconstrained matrix

$$\mathbf{a} = \begin{pmatrix} a_{1,1} & a_{1,2} & \dots & a_{1,K} \\ a_{2,1} & a_{2,2} & \dots & a_{2,K} \\ \vdots & \vdots & \ddots & \vdots \\ a_{N,1} & a_{N,2} & \dots & a_{N,K} \end{pmatrix}$$

and define

$$\phi = \text{softmax}(\mathbf{a})$$

where the softmax acts row-wise. Specifically:

$$\phi_{i,j} = \frac{\exp(a_{i,j})}{\sum_{j=1}^K \exp(a_{i,j})}$$

The exact **E-step**

$$\gamma_{i,k} = p(z^{(i)} | \mathbf{x}^{(i)}, \theta)$$

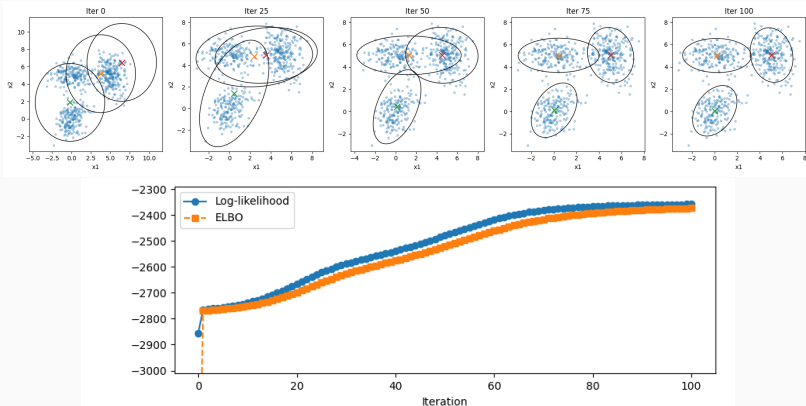
is replaced by a few gradient ascent steps:

$$\mathbf{a} \leftarrow \mathbf{a} + \eta \nabla_{\mathbf{a}} \text{ELBO}(\theta, \phi(\mathbf{a}))$$

This is fairly easy to implement, either by hand or using autodiff.

The **M-step** is exactly as before, but using $\phi_{i,j}$ instead of $\gamma_{i,j}$

- $\mu_k \leftarrow \frac{\sum_{i=1}^N \phi_{i,k} \mathbf{x}^{(i)}}{\sum_{i=1}^N \phi_{i,k}}$
- $\Sigma_k \leftarrow \frac{\sum_{i=1}^N \phi_{i,k} (\mathbf{x}^{(i)} - \mu_k)(\mathbf{x}^{(i)} - \mu_k)^T}{\sum_{i=1}^N \phi_{i,k}}$
- $w_k \leftarrow \frac{\sum_{i=1}^N \phi_{i,k}}{N}$



Here 5 gradient ascend steps are done in each EM iteration. Convergence is slower than for exact EM but achieves a similar solution (compare with last lecture).