

# Deep Generative Models

## Lecture 3

Roman Isachenko



Ozon Masters

Spring, 2021

# Recap of previous lecture

## MLE problem

$$\theta^* = \arg \max_{\theta} p(\mathbf{X}|\theta) = \arg \max_{\theta} \prod_{i=1}^n p(\mathbf{x}_i|\theta) = \arg \max_{\theta} \sum_{i=1}^n \log p(\mathbf{x}_i|\theta).$$

## Challenge

$p(\mathbf{x}|\theta)$  could be intractable.

## LVM

Introduce latent variable  $\mathbf{z}$  for each sample  $\mathbf{x}$

$$p(\mathbf{x}, \mathbf{z}|\theta) = p(\mathbf{x}|\mathbf{z}, \theta)p(\mathbf{z}); \quad \log p(\mathbf{x}, \mathbf{z}|\theta) = \log p(\mathbf{x}|\mathbf{z}, \theta) + \log p(\mathbf{z}).$$

$$p(\mathbf{x}|\theta) = \int p(\mathbf{x}, \mathbf{z}|\theta) d\mathbf{z} = \int p(\mathbf{x}|\mathbf{z}, \theta)p(\mathbf{z}) d\mathbf{z}.$$

## Motivation

The distributions  $p(\mathbf{x}|\mathbf{z}, \theta)$  and  $p(\mathbf{z})$  could be quite simple.

# Recap of previous lecture

## Incomplete likelihood maximization

$$\theta^* = \arg \max_{\theta} \log p(\mathbf{X}|\theta) = \arg \max_{\theta} \log \sum_{i=1}^n \int p(\mathbf{x}_i|\mathbf{z}_i, \theta) p(\mathbf{z}_i) d\mathbf{z}_i.$$

## Variational lower bound

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

## Evidence Lower Bound (ELBO)

$$\mathcal{L}(q, \theta) = \mathbb{E}_q \log p(\mathbf{x}|\mathbf{z}, \theta) - KL(q(\mathbf{z})||p(\mathbf{z}))$$

Instead of maximizing incomplete likelihood, maximize ELBO  
(equivalently minimize KL)

$$\max_{\theta} p(\mathbf{x}|\theta) \quad \rightarrow \quad \max_{q, \theta} \mathcal{L}(q, \theta) \equiv \min_{q, \theta} KL(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x}, \theta)).$$

# Recap of previous lecture

## EM algorithm (block-coordinate optimization)

- ▶ Initialize  $\theta^*$ ;

- ▶ E-step

$$q(\mathbf{z}) = \arg \max_q \mathcal{L}(q, \theta^*) = \arg \min_q KL(q||p) = p(\mathbf{z}|\mathbf{x}, \theta^*);$$

- ▶  $p(\mathbf{z}|\mathbf{x}, \theta^*)$  could be **intractable**;
- ▶  $q(\mathbf{z})$  is different for each object  $\mathbf{x}$ .

- ▶ M-step

$$\theta^* = \arg \max_{\theta} \mathcal{L}(q, \theta);$$

- ▶ Repeat E-step and M-step until convergence.

## Amortized variational inference

Restrict a family of all possible distributions  $q(\mathbf{z})$  to a particular parametric class  $q(\mathbf{z}|\mathbf{x}, \phi)$  conditioned on samples  $\mathbf{x}$  with parameters  $\phi$ .

# Variational EM-algorithm

## ELBO

$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \mathcal{L}(q, \boldsymbol{\theta}) + KL(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})) \geq \mathcal{L}(q, \boldsymbol{\theta}).$$

### ► E-step

$$\phi_k = \phi_{k-1} + \eta \nabla_{\phi} \mathcal{L}(\phi, \boldsymbol{\theta}_{k-1})|_{\phi=\phi_{k-1}},$$

where  $\phi$  – parameters of variational distribution  $q(\mathbf{z}|\mathbf{x}, \phi)$ .

### ► M-step

$$\boldsymbol{\theta}_k = \boldsymbol{\theta}_{k-1} + \eta \nabla_{\boldsymbol{\theta}} \mathcal{L}(\phi_k, \boldsymbol{\theta})|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{k-1}},$$

where  $\boldsymbol{\theta}$  – parameters of the generative distribution  $p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta})$ .

Now all we have to do is to obtain two gradients  $\nabla_{\phi} \mathcal{L}(\phi, \boldsymbol{\theta})$ ,  $\nabla_{\boldsymbol{\theta}} \mathcal{L}(\phi, \boldsymbol{\theta})$ .

## Challenge

Number of samples  $n$  could be huge (we need to derive unbiased stochastic gradients).

## ELBO gradient (M-step, $\nabla_{\theta} \mathcal{L}(\phi, \theta)$ )

$$\sum_{i=1}^n \mathcal{L}_i(\phi, \theta) = \sum_{i=1}^n \mathbb{E}_q [\log p(\mathbf{x}_i, \mathbf{z}_i | \theta) - \log q(\mathbf{z}_i | \mathbf{x}_i, \phi)] \rightarrow \max_{\phi, \theta}.$$

Optimization w.r.t.  $\theta$ : **mini-batching** (1) + **Monte-Carlo** estimation (2)

$$\begin{aligned} \nabla_{\theta} \sum_{i=1}^n \mathcal{L}_i(\phi, \theta) &= \sum_{i=1}^n \int q(\mathbf{z}_i | \mathbf{x}_i, \phi) \nabla_{\theta} \log p(\mathbf{x}_i | \mathbf{z}_i, \theta) d\mathbf{z}_i \\ &\stackrel{(1)}{\approx} n \int q(\mathbf{z}_i | \mathbf{x}_i, \phi) \nabla_{\theta} \log p(\mathbf{x}_i | \mathbf{z}_i, \theta) d\mathbf{z}_i, \quad i \sim U[1, n] \\ &\stackrel{(2)}{\approx} n \nabla_{\theta} \log p(\mathbf{x}_i | \mathbf{z}_i^*, \theta), \quad \mathbf{z}_i^* \sim q(\mathbf{z}_i | \mathbf{x}_i, \phi). \end{aligned}$$

**Monte-Carlo** estimation (2):

$$\mathbb{E}_q f(\mathbf{z}) = \int q(\mathbf{z}) f(\mathbf{z}) d\mathbf{z} \approx f(\mathbf{z}^*), \text{ where } \mathbf{z}^* \sim q(\mathbf{z}).$$

## ELBO gradient (E-step, $\nabla_{\phi} \mathcal{L}(\phi, \theta)$ )

$$\sum_{i=1}^n \mathcal{L}_i(\phi, \theta) = \sum_{i=1}^n \mathbb{E}_q [\log p(\mathbf{x}_i, \mathbf{z}_i | \theta) - \log q(\mathbf{z}_i | \mathbf{x}_i, \phi)] \rightarrow \max_{\phi, \theta}.$$

### Challenge

Difference from M-step: density function  $q(\mathbf{z} | \mathbf{x}, \phi)$  depends on the parameters  $\phi$ , it is impossible to use the Monte-Carlo estimation:

$$\begin{aligned} \nabla_{\phi} \mathcal{L}(\phi, \theta) &= \nabla_{\phi} \int q(\mathbf{z} | \mathbf{x}, \phi) [\log p(\mathbf{x}_i, \mathbf{z}_i | \theta) - \log q(\mathbf{z}_i | \mathbf{x}_i, \phi)] d\mathbf{z} \\ &\neq \int q(\mathbf{z} | \mathbf{x}, \phi) \nabla_{\phi} [\log p(\mathbf{x}_i, \mathbf{z}_i | \theta) - \log q(\mathbf{z}_i | \mathbf{x}_i, \phi)] d\mathbf{z} \end{aligned}$$

### Solution

Reparametrization trick for  $q(\mathbf{z} | \mathbf{x}, \phi)$  to allow the expectation is independent of parameters  $\phi$ .

## ELBO gradient (E-step, $\nabla_{\phi} \mathcal{L}(\phi, \theta)$ )

Reparametrization trick

$$f(\xi) = \int q(\eta|\xi) h(\eta) d\eta$$

Let  $\eta = g(\xi, \epsilon)$ , where  $g$  is a deterministic function,  $\epsilon$  is a random variable with a density function  $r(\epsilon)$ .

$$f(\xi) = \int q(\eta|\xi) h(\eta) d\eta = \int r(\epsilon) h(g(\xi, \epsilon)) d\epsilon \approx h(g(\xi, \epsilon^*)), \quad \epsilon^* \sim r(\epsilon).$$

$$\nabla_{\xi} f(\xi) \approx \nabla_{\xi} h(g(\xi, \epsilon^*)), \quad \epsilon^* \sim r(\epsilon).$$

### Example

$$q(\eta|\xi) = \mathcal{N}(\eta|\mu, \sigma^2), \quad r(\epsilon) = \mathcal{N}(\epsilon|0, 1), \quad \eta = \sigma \cdot \epsilon + \mu, \quad \xi = [\mu, \sigma].$$

$$\epsilon \sim r(\epsilon), \quad \mathbf{z} = g(\mathbf{x}, \epsilon, \phi), \quad \mathbf{z} \sim q(\mathbf{z}|\mathbf{x}, \phi).$$

$$\nabla_{\phi} \sum_{i=1}^n \mathcal{L}_i(\phi, \theta) = \sum_{i=1}^n \nabla_{\phi} \int q(\mathbf{z}_i|\mathbf{x}_i, \phi) \log p(\mathbf{x}_i|\mathbf{z}_i, \theta) d\mathbf{z}_i - \nabla_{\phi} KL$$



## ELBO gradient (E-step, $\nabla_{\phi} \mathcal{L}(\phi, \theta)$ )

$$\begin{aligned}\nabla_{\phi} \sum_{i=1}^n \mathcal{L}_i(\phi, \theta) &= \sum_{i=1}^n \nabla_{\phi} \int q(\mathbf{z}_i | \mathbf{x}_i, \phi) \log p(\mathbf{x}_i | \mathbf{z}_i, \theta) d\mathbf{z}_i - \nabla_{\phi} KL \\ &\approx n \nabla_{\phi} \int r(\epsilon) \log p(\mathbf{x}_i | g(\mathbf{x}_i, \epsilon, \phi), \theta) d\epsilon - \nabla_{\phi} KL, \quad i \sim U[1, n] \\ &\approx n \nabla_{\phi} \log p(\mathbf{x}_i | g(\mathbf{x}_i, \epsilon^*, \phi), \theta) - \nabla_{\phi} KL, \quad \epsilon^* \sim r(\epsilon).\end{aligned}$$

### Variational assumption

$$q(\mathbf{z} | \mathbf{x}, \phi) = \mathcal{N}(\mu(\mathbf{x}), \sigma(\mathbf{x})).$$

$$\mathbf{z} = g(\mathbf{x}, \epsilon, \phi) = \sigma(\mathbf{x}) \cdot \epsilon + \mu(\mathbf{x}).$$

$\nabla_{\phi} KL(q(\mathbf{z} | \mathbf{x}, \phi) || p(\mathbf{z}))$  has an analytical solution.

# Variational autoencoder (VAE)

## Final algorithm

- ▶ pick  $i \sim U[1, n]$ ;
- ▶ compute a stochastic gradient w.r.t.  $\phi$

$$\nabla_{\phi} \mathcal{L}(\phi, \theta) = n \nabla_{\phi} \log p(\mathbf{x}_i | g(\mathbf{x}_i, \epsilon^*, \phi), \theta) - \nabla_{\phi} KL(q(\mathbf{z}_i | \mathbf{x}_i, \phi) || p(\mathbf{z}_i)), \quad \epsilon^* \sim r(\epsilon);$$

- ▶ compute a stochastic gradient w.r.t.  $\theta$

$$\nabla_{\theta} \mathcal{L}(\phi, \theta) = n \nabla_{\theta} \log p(\mathbf{x}_i | \mathbf{z}_i^*, \theta), \quad \mathbf{z}_i^* \sim q(\mathbf{z}_i | \mathbf{x}_i, \phi);$$

- ▶ update  $\theta, \phi$  according to the selected optimization method (SGD, Adam, RMSProp).

# Variational autoencoder (VAE)

- ▶ Encoder  $q(\mathbf{z}|\mathbf{x}, \phi) = \text{NN}_e(\mathbf{x}, \phi)$  outputs  $\mu(\mathbf{x})$  and  $\sigma(\mathbf{x})$ .
- ▶ Decoder  $p(\mathbf{x}|\mathbf{z}, \theta) = \text{NN}_d(\mathbf{z}, \theta)$  outputs parameters of the sample distribution.

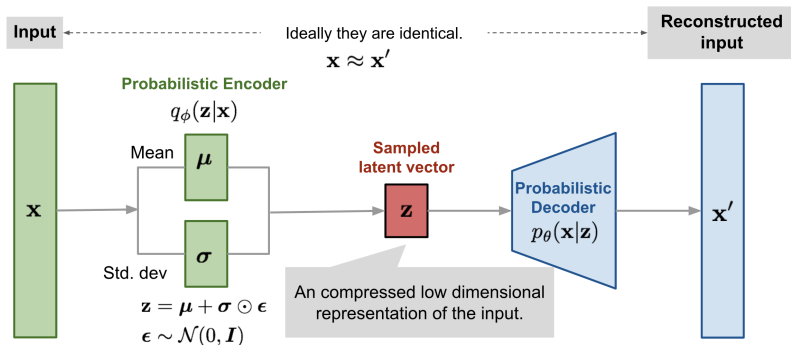
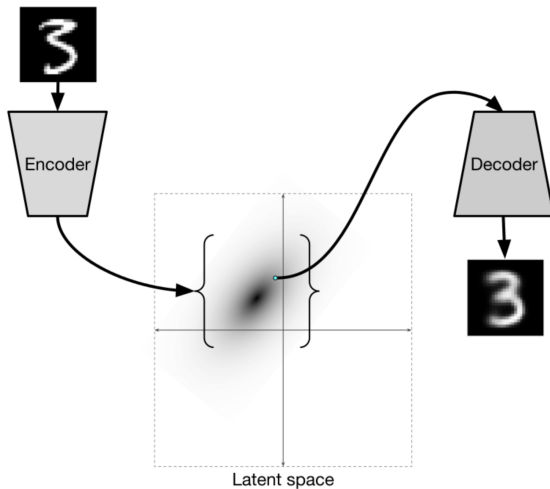


image credit:

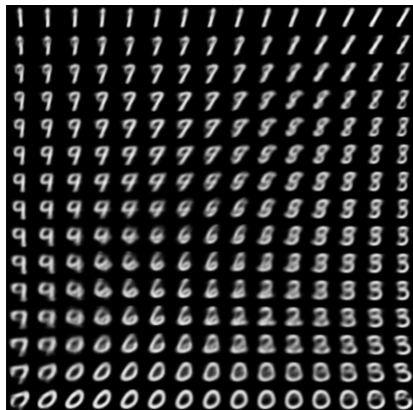
<https://lilianweng.github.io/lil-log/2018/08/12/from-autoencoder-to-beta-vae.html>

# Variational Autoencoder



# Variational Autoencoder

Generation objects by sampling the latent space  $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$



# Bayesian framework

## Bayes theorem

$$p(\mathbf{t}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{t})p(\mathbf{t})}{p(\mathbf{x})} = \frac{p(\mathbf{x}|\mathbf{t})p(\mathbf{t})}{\int p(\mathbf{x}|\mathbf{t})p(\mathbf{t})d\mathbf{t}}$$

- ▶  $\mathbf{x}$  – observed variables;
- ▶  $\mathbf{t}$  – unobserved variables (latent variables/parameters);
- ▶  $p(\mathbf{x}|\mathbf{t})$  – likelihood;
- ▶  $p(\mathbf{x})$  – evidence;
- ▶  $p(\mathbf{t})$  – prior;
- ▶  $p(\mathbf{t}|\mathbf{x})$  – posterior.

# Variational Lower Bound

We are given the set of objects  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n$ . The goal is to perform bayesian inference on the latent variables  $\mathbf{T} = \{\mathbf{t}_i\}_{i=1}^n$ .

## Evidence Lower Bound (ELBO)

$$\begin{aligned}\log p(\mathbf{X}) &= \log \frac{p(\mathbf{X}, \mathbf{T})}{p(\mathbf{T}|\mathbf{X})} = \\ &= \int q(\mathbf{T}) \log \frac{p(\mathbf{X}, \mathbf{T})}{p(\mathbf{T}|\mathbf{X})} d\mathbf{T} = \int q(\mathbf{T}) \log \frac{p(\mathbf{X}, \mathbf{T})q(\mathbf{T})}{p(\mathbf{T}|\mathbf{X})q(\mathbf{T})} d\mathbf{T} = \\ &= \int q(\mathbf{T}) \log \frac{p(\mathbf{X}, \mathbf{T})}{q(\mathbf{T})} d\mathbf{T} + \int q(\mathbf{T}) \log \frac{q(\mathbf{T})}{p(\mathbf{T}|\mathbf{X})} d\mathbf{T} = \\ &= \mathcal{L}(q) + KL(q(\mathbf{T})||p(\mathbf{T}|\mathbf{X})) \geq \mathcal{L}(q).\end{aligned}$$

We would like to maximize lower bound  $\mathcal{L}(q)$ .

# Mean field approximation

## Independence assumption

$$q(\mathbf{T}) = \prod_{i=1}^k q_i(\mathbf{T}_i), \quad \mathbf{T} = [\mathbf{T}_1, \dots, \mathbf{T}_k], \quad \mathbf{T}_j = \{\mathbf{t}_{ij}\}_{i=1}^n, \quad \mathbf{t}_i = \{\mathbf{T}_{ij}\}_{j=1}^k.$$

## Block coordinate optimization of ELBO for $q_j(\mathbf{T}_j)$

$$\begin{aligned} \mathcal{L}(q) &= \int q(\mathbf{T}) \log \frac{p(\mathbf{X}, \mathbf{T})}{q(\mathbf{T})} d\mathbf{T} = \int \prod_{i=1}^k q_i(\mathbf{T}_i) \log \frac{p(\mathbf{X}, \mathbf{T})}{\prod_{i=1}^k q_i(\mathbf{T}_i)} \prod_{i=1}^k d\mathbf{T}_i = \\ &= \int \prod_{i=1}^k q_i \log p(\mathbf{X}, \mathbf{T}) \prod_{i=1}^k d\mathbf{T}_i - \sum_{i=1}^k \int \prod_{j=1}^k q_j \log q_i \prod_{i=1}^k d\mathbf{T}_i = \\ &= \int q_j \left[ \int \log p(\mathbf{X}, \mathbf{T}) \prod_{i \neq j} q_i d\mathbf{T}_i \right] d\mathbf{T}_j - \\ &\quad - \int q_j \log q_j d\mathbf{T}_j + \text{const}(q_j) \rightarrow \max_{q_j} \end{aligned}$$



# Mean field approximation

## Block coordinate optimization of ELBO for $q_j(\mathbf{T}_j)$

$$\begin{aligned}\mathcal{L}(q) &= \int q_j \left[ \int \log p(\mathbf{X}, \mathbf{T}) \prod_{i \neq j} q_i d\mathbf{T}_i \right] d\mathbf{T}_j - \int q_j \log q_j d\mathbf{T}_j + \text{const}(q_j) = \\ &= \int q_j \log \hat{p}(\mathbf{X}, \mathbf{T}_j) d\mathbf{T}_j - \int q_j \log q_j d\mathbf{T}_j + \text{const}(q_j) \rightarrow \max_{q_j},\end{aligned}$$

$$\text{where } \log \hat{p}(\mathbf{X}, \mathbf{T}_j) = \mathbb{E}_{i \neq j} \log p(\mathbf{X}, \mathbf{T}) + \text{const}(q_j)$$

$$\mathbb{E}_{i \neq j} \log p(\mathbf{X}, \mathbf{T}) = \int \log p(\mathbf{X}, \mathbf{T}) \prod_{i \neq j} q_i d\mathbf{T}_i.$$

$$\begin{aligned}\mathcal{L}(q) &= \int q_j(\mathbf{T}_j) \log \hat{p}(\mathbf{X}, \mathbf{T}_j) d\mathbf{T}_j - \int q_j(\mathbf{T}_j) \log q_j(\mathbf{T}_j) d\mathbf{T}_j + \text{const}(q_j) = \\ &= \int q_j(\mathbf{T}_j) \log \frac{\hat{p}(\mathbf{X}, \mathbf{T}_j)}{q_j(\mathbf{T}_j)} d\mathbf{T}_j + \text{const}(q_j) = \\ &= -KL(q_j(\mathbf{T}_j) || \hat{p}(\mathbf{X}, \mathbf{T}_j)) + \text{const}(q_j) \rightarrow \max_{q_j}.\end{aligned}$$

# Mean field approximation

## Independence assumption

$$q(\mathbf{T}) = \prod_{i=1}^k q_i(\mathbf{T}_i), \quad \mathbf{T} = [\mathbf{T}_1, \dots, \mathbf{T}_k], \quad \mathbf{T}_j = \{\mathbf{t}_{ij}\}_{i=1}^n.$$

## ELBO

$$\mathcal{L}(q) = -KL(q_j(\mathbf{T}_j) \parallel \hat{p}(\mathbf{X}, \mathbf{T}_j)) + \text{const}(q_j) \rightarrow \max_{q_j}.$$

## Solution

$$q_j(\mathbf{T}_j) = \hat{p}(\mathbf{X}, \mathbf{T}_j)$$

$$\log \hat{p}(\mathbf{X}, \mathbf{T}_j) = \mathbb{E}_{i \neq j} \log p(\mathbf{X}, \mathbf{T}) + \text{const}$$

$$\log q_j(\mathbf{T}_j) = \mathbb{E}_{i \neq j} \log p(\mathbf{X}, \mathbf{T}) + \text{const}$$

# Mean field approximation

## ELBO

$$\mathcal{L}(q) = -KL(q_j(\mathbf{T}_j) || \hat{p}(\mathbf{X}, \mathbf{T}_j)) + \text{const}(q_j) \rightarrow \max_{q_j}.$$

## Solution

$$\log q_j(\mathbf{T}_j) = \mathbb{E}_{i \neq j} \log p(\mathbf{X}, \mathbf{T}) + \text{const}$$

Let assume the following factorization:  $\mathbf{T} = [\mathbf{T}_1, \mathbf{T}_2] = [\mathbf{Z}, \boldsymbol{\theta}]$ , and restrict the class of probability distribution for  $\boldsymbol{\theta}$  to Dirac delta functions:

$$q_2 = q(\mathbf{T}_2) = q(\boldsymbol{\theta}) = \delta(\boldsymbol{\theta} - \boldsymbol{\theta}_0).$$

Under the restrictions the exact solution for  $q_2$  is not reached.

# Mean field approximation

## General solution

$$\log q_j(\mathbf{T}_j) = \mathbb{E}_{i \neq j} \log p(\mathbf{X}, \mathbf{T}) + \text{const}$$

## Solution for $q_1 = q(\mathbf{Z})$

$$\begin{aligned} \log q(\mathbf{Z}) &= \int q(\boldsymbol{\theta}) \log p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\theta}) d\boldsymbol{\theta} + \text{const} = \\ &= \int \delta(\boldsymbol{\theta} - \boldsymbol{\theta}_0) \log p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\theta}) d\boldsymbol{\theta} + \text{const} = \\ &= \log p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}_0) + \text{const}. \end{aligned}$$

## EM-algorithm (E-step)

$$q(\mathbf{Z}) = \arg \max_q \mathcal{L}(q, \boldsymbol{\theta}^*) = \arg \min_q KL(q||p) = p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^*).$$

# Mean field approximation

## ELBO

$$\mathcal{L}(q) = -KL(q_j(\mathbf{T}_j) || \hat{p}(\mathbf{X}, \mathbf{T}_j)) + \text{const}(q_j) \rightarrow \max_{q_j}.$$

ELBO maximization w.r.t.  $q_2 \equiv \theta_0$

$$\begin{aligned}\mathcal{L}(q_2) &= -KL(q(\boldsymbol{\theta}) || \hat{p}(\mathbf{X}, \boldsymbol{\theta})) + \text{const}(\boldsymbol{\theta}_0) \\&= \int q(\boldsymbol{\theta}) \log \frac{\hat{p}(\mathbf{X}, \boldsymbol{\theta})}{q(\boldsymbol{\theta})} d\boldsymbol{\theta} + \text{const}(\boldsymbol{\theta}_0) \\&= \int q(\boldsymbol{\theta}) \log \hat{p}(\mathbf{X}, \boldsymbol{\theta}) d\boldsymbol{\theta} - \int q(\boldsymbol{\theta}) \log q(\boldsymbol{\theta}) d\boldsymbol{\theta} + \text{const}(\boldsymbol{\theta}_0) \\&= \int \delta(\boldsymbol{\theta} - \boldsymbol{\theta}_0) \log \hat{p}(\mathbf{X}, \boldsymbol{\theta}) d\boldsymbol{\theta} - \int \delta \log \delta d\boldsymbol{\theta} + \text{const}(\boldsymbol{\theta}_0) \\&= \int \delta(\boldsymbol{\theta} - \boldsymbol{\theta}_0) \log \hat{p}(\mathbf{X}, \boldsymbol{\theta}) d\boldsymbol{\theta} + \text{const}(\boldsymbol{\theta}_0)\end{aligned}$$

# Mean field approximation

ELBO maximization w.r.t.  $q_2 \equiv \theta_0$

$$\mathcal{L}(q_2) = \int \delta(\boldsymbol{\theta} - \boldsymbol{\theta}_0) \log \hat{p}(\mathbf{X}, \boldsymbol{\theta}) d\boldsymbol{\theta} + \text{const}(\boldsymbol{\theta}_0) = \log \hat{p}(\mathbf{X}, \boldsymbol{\theta}_0).$$

$$\log \hat{p}(\mathbf{X}, \mathbf{T}_j) = \mathbb{E}_{i \neq j} \log p(\mathbf{X}, \mathbf{T}) + \text{const}$$

$$\begin{aligned} \log \hat{p}(\mathbf{X}, \boldsymbol{\theta}) &= \mathbb{E}_{q_1} \log p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\theta}) + \text{const} \\ &= \int q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}) d\mathbf{Z} + \log p(\boldsymbol{\theta}) + \text{const} \end{aligned}$$

EM-algorithm (M-step)

$$\mathcal{L}(q, \boldsymbol{\theta}) = \int q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta})}{q(\mathbf{Z})} d\mathbf{Z} \rightarrow \max_{\boldsymbol{\theta}}$$

# Mean field approximation

## Solution

$$\log q_j(\mathbf{T}_j) = \mathbb{E}_{i \neq j} \log p(\mathbf{X}, \mathbf{T}) + \text{const}$$

## EM algorithm (special case)

- ▶ Initialize  $\theta^*$ ;
- ▶ E-step

$$q(\mathbf{Z}) = \arg \max_q \mathcal{L}(q, \theta^*) = \arg \min_q KL(q||p) = p(\mathbf{Z}|\mathbf{X}, \theta^*);$$

- ▶ M-step

$$\theta^* = \arg \max_{\theta} \mathcal{L}(q, \theta);$$

- ▶ Repeat E-step and M-step until convergence.

# Summary

- ▶ Bayesian inference is a generalization of most common machine learning tasks. It allows to construct MLE, MAP and bayesian inference, to compare models complexity and many-many more cool stuff.
- ▶ LVM introduce latent representation of observed samples to make model more interpretable.
- ▶ LVM maximizes variational evidence lower bound to find MLE of model parameters.
- ▶ ELBO maximization is performed by the general variational EM algorithm.
- ▶ Amortized inference allows to efficiently compute stochastic gradients for ELBO and to use deep neural networks for  $q(\mathbf{z}|\mathbf{x}, \phi)$  and  $p(\mathbf{x}|\mathbf{z}, \theta)$ .
- ▶ The VAE model is an LVM with an encoder network for  $q(\mathbf{z}|\mathbf{x}, \phi)$  and a decoder network for  $p(\mathbf{x}|\mathbf{z}, \theta)$ .



# Summary

- ▶ Latent variable models introduce latent variables to the initial probabilistic model to make distribution  $p(\mathbf{x}|\boldsymbol{\theta})$  tractable.
- ▶ To solve the MLE problem LVM optimizes the variational lower bound.
- ▶ The EM-algorithm is an iterative algorithm which allows to optimize the variational lower bound.
- ▶ VAE model is an LVM, where the encoder gives the variational distribution, the decoder defines the likelihood model.
- ▶ The mean field approximation is a general form of variational inference (the EM-algorithm is just a special case).