# Deep Generative Models

## Lecture 3

Roman Isachenko

**Moscow Institute of Physics and Technology**
**Yandex School of Data Analysis**

2025, Autumn

# Recap of Previous Lecture

### Jacobian Matrix

Given a differentiable function $\mathbf{f} : \mathbb{R}^m \to \mathbb{R}^m$,

$$\mathbf{z} = \mathbf{f}(\mathbf{x}), \quad \mathbf{J} = \frac{\partial \mathbf{z}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial z_1}{\partial x_1} & \cdots & \frac{\partial z_1}{\partial x_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial z_m}{\partial x_1} & \cdots & \frac{\partial z_m}{\partial x_m} \end{pmatrix} \in \mathbb{R}^{m \times m}$$
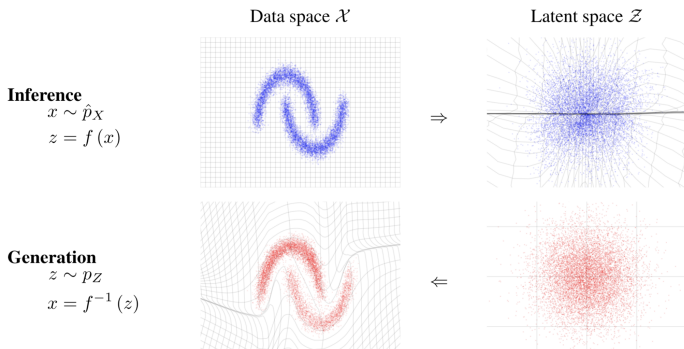
### Change of Variables Theorem (CoV)

Let $\mathbf{x}$ be a random variable with density $p(\mathbf{x})$, and $\mathbf{f} : \mathbb{R}^m \to \mathbb{R}^m$ a differentiable invertible mapping. If $\mathbf{z} = \mathbf{f}(\mathbf{x})$ and $\mathbf{x} = \mathbf{f}^{-1}(\mathbf{z}) = \mathbf{g}(\mathbf{z})$, then

$$p(\mathbf{x}) = p(\mathbf{z})|\det(\mathbf{J_f})| = p(\mathbf{z}) \left| \det\left( \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \right) \right| = p(\mathbf{f}(\mathbf{x})) \left| \det\left( \frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} \right) \right|$$

$$p(\mathbf{z}) = p(\mathbf{x})|\det(\mathbf{J_g})| = p(\mathbf{x}) \left| \det\left( \frac{\partial \mathbf{x}}{\partial \mathbf{z}} \right) \right| = p(\mathbf{g}(\mathbf{z})) \left| \det\left( \frac{\partial \mathbf{g}(\mathbf{z})}{\partial \mathbf{z}} \right) \right|.$$

---

Dinh L., Sohl-Dickstein J., Bengio S. Density Estimation Using Real NVP, 2016

# Recap of Previous Lecture

### Definition
A normalizing flow is a *differentiable*, *invertible* transformation that maps data **x** to noise **z**.



Data space $\mathcal{X}$      Latent space $\mathcal{Z}$

**Inference**
$x \sim \hat{p}_X$
$z = f(x)$

$\Rightarrow$

**Generation**
$z \sim p_Z$
$x = f^{-1}(z)$

$\Leftarrow$

### Log-Likelihood

$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \log p(\mathbf{f}_K \circ \cdots \circ \mathbf{f}_1(\mathbf{x})) + \sum_{k=1}^{K} \log |\det(\mathbf{J}_{\mathbf{f}_k})|$$

*Dinh L., Sohl-Dickstein J., Bengio S. Density Estimation Using Real NVP, 2016*

# Recap of Previous Lecture

### Flow Log-Likelihood

$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \log p(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})) + \log|\det(\mathbf{J_f})|$$

One significant challenge is efficiently computing the Jacobian determinant.

### Linear Flows

$$\mathbf{z} = \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}) = \mathbf{W}\mathbf{x}, \quad \mathbf{W} \in \mathbb{R}^{m \times m}, \quad \boldsymbol{\theta} = \mathbf{W}, \quad \mathbf{J_f} = \mathbf{W}^T$$

▶ LU Decomposition:

$$\mathbf{W} = \mathbf{PLU}.$$

▶ QR Decomposition:

$$\mathbf{W} = \mathbf{QR}.$$

Decomposition is performed only once during initialization. Then the decomposed matrices ($\mathbf{P}, \mathbf{L}, \mathbf{U}$ or $\mathbf{Q}, \mathbf{R}$) are optimized.

Kingma D. P., et al. Glow: Generative Flow With Invertible 1x1 Convolutions, 2018
Hoogeboom E., et al. Emerging Convolutions for Generative Normalizing Flows, 2019

# Recap of Previous Lecture

Consider an autoregressive model:

$$p(\mathbf{x}|\boldsymbol{\theta}) = \prod_{j=1}^{m} p(x_j|\mathbf{x}_{1:j-1}, \boldsymbol{\theta}), \quad p(x_j|\mathbf{x}_{1:j-1}, \boldsymbol{\theta}) = \mathcal{N}\left(\mu_{j,\boldsymbol{\theta}}(\mathbf{x}_{1:j-1}), \sigma_{j,\boldsymbol{\theta}}^2(\mathbf{x}_{1:j-1})\right).$$

Gaussian Autoregressive Normalizing Flow

$$\mathbf{x} = \mathbf{g}_{\boldsymbol{\theta}}(\mathbf{z}) \quad \Rightarrow \quad x_j = \sigma_{j,\boldsymbol{\theta}}(\mathbf{x}_{1:j-1}) \cdot z_j + \mu_{j,\boldsymbol{\theta}}(\mathbf{x}_{1:j-1}).$$

$$\mathbf{z} = \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}) \quad \Rightarrow \quad z_j = (x_j - \mu_{j,\boldsymbol{\theta}}(\mathbf{x}_{1:j-1})) \cdot \frac{1}{\sigma_{j,\boldsymbol{\theta}}(\mathbf{x}_{1:j-1})}.$$

▶ This transformation is both **invertible** and **differentiable**, mapping $p(\mathbf{z})$ to $p(\mathbf{x}|\boldsymbol{\theta})$.

▶ The Jacobian matrix for this transformation is triangular.

The generative function $\mathbf{g}_{\boldsymbol{\theta}}(\mathbf{z})$ is **sequential**, while the inference function $\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})$ is **not sequential**.

Papamakarios G., Pavlakou T., Murray I. Masked Autoregressive Flow for Density Estimation, 2017

# Recap of Previous Lecture

Let us partition $\mathbf{x}$ and $\mathbf{z}$ into two groups:

$$\mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2] = [\mathbf{x}_{1:d}, \mathbf{x}_{d+1:m}]; \quad \mathbf{z} = [\mathbf{z}_1, \mathbf{z}_2] = [\mathbf{z}_{1:d}, \mathbf{z}_{d+1:m}].$$

## Coupling Layer

$$\begin{cases} \mathbf{x}_1 = \mathbf{z}_1; \\ \mathbf{x}_2 = \mathbf{z}_2 \odot \boldsymbol{\sigma}_{\boldsymbol{\theta}}(\mathbf{z}_1) + \boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{z}_1). \end{cases} \qquad \begin{cases} \mathbf{z}_1 = \mathbf{x}_1; \\ \mathbf{z}_2 = (\mathbf{x}_2 - \boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{x}_1)) \odot \frac{1}{\boldsymbol{\sigma}_{\boldsymbol{\theta}}(\mathbf{x}_1)}. \end{cases}$$

Both density estimation and sampling require just a single pass!

## Jacobian

$$\det \left( \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \right) = \det \begin{pmatrix} \mathbf{I}_d & 0_{d \times (m-d)} \\ \frac{\partial \mathbf{z}_2}{\partial \mathbf{x}_1} & \frac{\partial \mathbf{z}_2}{\partial \mathbf{x}_2} \end{pmatrix} = \prod_{j=1}^{m-d} \frac{1}{\sigma_{j,\theta}(\mathbf{x}_1)}.$$

A coupling layer is a special instance of an gaussian autoregressive normalizing flow.

*Dinh L., Sohl-Dickstein J., Bengio S. Density Estimation Using Real NVP, 2016*

# Outline

1. Forward and Reverse KL for NF

2. Latent Variable Models (LVM)

3. Variational Evidence Lower Bound (ELBO)

4. EM-Algorithm

# Outline

# Forward KL vs. Reverse KL

## Forward KL ($\equiv$ Maximum Likelihood Estimation)

$$\mathrm{KL}(\pi\|p) = \int \pi(\mathbf{x}) \log \frac{\pi(\mathbf{x})}{p(\mathbf{x}|\boldsymbol{\theta})} d\mathbf{x}$$
$$= -\mathbb{E}_{\pi(\mathbf{x})} \log p(\mathbf{x}|\boldsymbol{\theta}) + \mathrm{const} \to \min_{\boldsymbol{\theta}}$$

# Forward KL vs. Reverse KL

### Forward KL ($\equiv$ Maximum Likelihood Estimation)

$$\mathrm{KL}(\pi \| p) = \int \pi(\mathbf{x}) \log \frac{\pi(\mathbf{x})}{p(\mathbf{x}|\boldsymbol{\theta})} d\mathbf{x}$$

$$= -\mathbb{E}_{\pi(\mathbf{x})} \log p(\mathbf{x}|\boldsymbol{\theta}) + \mathrm{const} \to \min_{\boldsymbol{\theta}}$$

### Forward KL for Normalizing Flows

$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \log p(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})) + \log |\det(\mathbf{J}_{\mathbf{f}})|$$

# Forward KL vs. Reverse KL

### Forward KL ($\equiv$ Maximum Likelihood Estimation)

$$\mathrm{KL}(\pi\|p) = \int \pi(\mathbf{x}) \log \frac{\pi(\mathbf{x})}{p(\mathbf{x}|\boldsymbol{\theta})} d\mathbf{x}$$
$$= -\mathbb{E}_{\pi(\mathbf{x})} \log p(\mathbf{x}|\boldsymbol{\theta}) + \mathsf{const} \to \min_{\boldsymbol{\theta}}$$

### Forward KL for Normalizing Flows

$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \log p(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})) + \log |\det(\mathbf{J}_{\mathbf{f}})|$$

$$\mathrm{KL}(\pi\|p) = -\mathbb{E}_{\pi(\mathbf{x})} \Big[ \log p(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})) + \log |\det(\mathbf{J}_{\mathbf{f}})| \Big] + \mathsf{const}$$

# Forward KL vs. Reverse KL

### Forward KL ($\equiv$ Maximum Likelihood Estimation)

$$\mathrm{KL}(\pi \| p) = \int \pi(\mathbf{x}) \log \frac{\pi(\mathbf{x})}{p(\mathbf{x}|\boldsymbol{\theta})} d\mathbf{x}$$
$$= -\mathbb{E}_{\pi(\mathbf{x})} \log p(\mathbf{x}|\boldsymbol{\theta}) + \text{const} \to \min_{\boldsymbol{\theta}}$$

### Forward KL for Normalizing Flows

$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \log p(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})) + \log |\det(\mathbf{J_f})|$$

$$\mathrm{KL}(\pi \| p) = -\mathbb{E}_{\pi(\mathbf{x})} \Big[ \log p(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})) + \log |\det(\mathbf{J_f})| \Big] + \text{const}$$

- We need to compute $\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})$ and its Jacobian.
- Access to the density $p(\mathbf{z})$ is required.
- The inverse function $\mathbf{g}_{\boldsymbol{\theta}}(\mathbf{z}) = \mathbf{f}_{\boldsymbol{\theta}}^{-1}(\mathbf{z})$ is required only for sampling from the normalizing flow.

# Forward KL vs. Reverse KL

## Reverse KL

$$\mathrm{KL}(p\|\pi) = \int p(\mathbf{x}|\boldsymbol{\theta}) \log \frac{p(\mathbf{x}|\boldsymbol{\theta})}{\pi(\mathbf{x})} d\mathbf{x}$$
$$= \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta})} \left[ \log p(\mathbf{x}|\boldsymbol{\theta}) - \log \pi(\mathbf{x}) \right] \rightarrow \min_{\boldsymbol{\theta}}$$

# Forward KL vs. Reverse KL

### Reverse KL

$$\mathrm{KL}(p\|\pi) = \int p(\mathbf{x}|\boldsymbol{\theta}) \log \frac{p(\mathbf{x}|\boldsymbol{\theta})}{\pi(\mathbf{x})} d\mathbf{x}$$
$$= \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta})} \left[ \log p(\mathbf{x}|\boldsymbol{\theta}) - \log \pi(\mathbf{x}) \right] \to \min_{\boldsymbol{\theta}}$$

### Reverse KL for Normalizing Flows (LOTUS Trick)

$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \log p(\mathbf{z}) + \log |\det(\mathbf{J_f})| = \log p(\mathbf{z}) - \log |\det(\mathbf{J_g})|$$

$$\mathrm{KL}(p\|\pi) = \mathbb{E}_{p(\mathbf{z})} \left[ \log p(\mathbf{z}) - \log |\det(\mathbf{J_g})| - \log \pi(\mathbf{g}_{\boldsymbol{\theta}}(\mathbf{z})) \right]$$

# Forward KL vs. Reverse KL

### Reverse KL

$$\mathrm{KL}(p\|\pi) = \int p(\mathbf{x}|\boldsymbol{\theta}) \log \frac{p(\mathbf{x}|\boldsymbol{\theta})}{\pi(\mathbf{x})} d\mathbf{x}$$
$$= \mathbb{E}_{p(\mathbf{x}|\boldsymbol{\theta})} \left[\log p(\mathbf{x}|\boldsymbol{\theta}) - \log \pi(\mathbf{x})\right] \to \min_{\boldsymbol{\theta}}$$

### Reverse KL for Normalizing Flows (LOTUS Trick)

$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \log p(\mathbf{z}) + \log |\det(\mathbf{J_f})| = \log p(\mathbf{z}) - \log |\det(\mathbf{J_g})|$$

$$\mathrm{KL}(p\|\pi) = \mathbb{E}_{p(\mathbf{z})} \left[\log p(\mathbf{z}) - \log |\det(\mathbf{J_g})| - \log \pi(\mathbf{g}_{\boldsymbol{\theta}}(\mathbf{z}))\right]$$

▶ We need to compute $\mathbf{g}_{\boldsymbol{\theta}}(\mathbf{z})$ and its Jacobian.

▶ Sampling from $p(\mathbf{z})$ is required (though direct evaluation is not), along with evaluating $\pi(\mathbf{x})$.

▶ Evaluating $\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})$ is not required.

# Normalizing Flows KL Duality

### Theorem
Fitting the NF model $p(\mathbf{x}|\boldsymbol{\theta})$ to a target distribution $\pi(\mathbf{x})$ via the forward KL (MLE) is equivalent to fitting the induced distribution $p(\mathbf{z}|\boldsymbol{\theta})$ to the base distribution $p(\mathbf{z})$ via the reverse KL:

$$\arg\min_{\boldsymbol{\theta}} \mathrm{KL}(\pi(\mathbf{x})\|p(\mathbf{x}|\boldsymbol{\theta})) = \arg\min_{\boldsymbol{\theta}} \mathrm{KL}(p(\mathbf{z}|\boldsymbol{\theta})\|p(\mathbf{z})).$$

Papamakarios G. et al. *Normalizing Flows for Probabilistic Modeling and Inference*, 2019

# Normalizing Flows KL Duality

### Theorem

Fitting the NF model $p(\mathbf{x}|\boldsymbol{\theta})$ to a target distribution $\pi(\mathbf{x})$ via the forward KL (MLE) is equivalent to fitting the induced distribution $p(\mathbf{z}|\boldsymbol{\theta})$ to the base distribution $p(\mathbf{z})$ via the reverse KL:
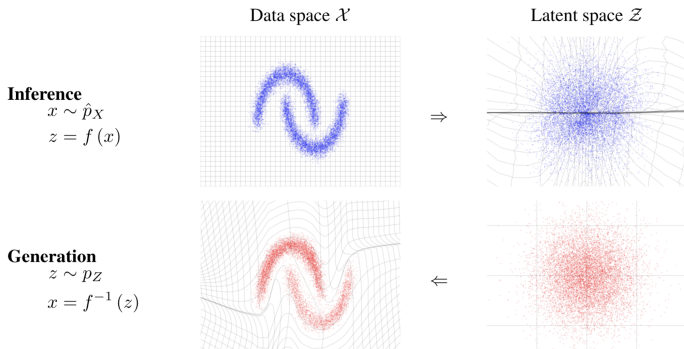
$$\arg\min_{\boldsymbol{\theta}} \mathrm{KL}(\pi(\mathbf{x})\|p(\mathbf{x}|\boldsymbol{\theta})) = \arg\min_{\boldsymbol{\theta}} \mathrm{KL}(p(\mathbf{z}|\boldsymbol{\theta})\|p(\mathbf{z})).$$



Data space $\mathcal{X}$            Latent space $\mathcal{Z}$

**Inference**
$x \sim \hat{p}_X$
$z = f(x)$

$\Rightarrow$

**Generation**
$z \sim p_Z$
$x = f^{-1}(z)$

$\Leftarrow$

Papamakarios G. et al. *Normalizing Flows for Probabilistic Modeling and Inference*, 2019

# Normalizing Flows KL Duality

## Theorem

$$\arg\min_{\boldsymbol{\theta}} \mathrm{KL}(\pi(\mathbf{x}) \| p(\mathbf{x}|\boldsymbol{\theta})) = \arg\min_{\boldsymbol{\theta}} \mathrm{KL}(p(\mathbf{z}|\boldsymbol{\theta}) \| p(\mathbf{z})).$$

## Proof

- $\mathbf{z} \sim p(\mathbf{z})$, $\mathbf{x} = \mathbf{g}_{\theta}(\mathbf{z})$, so $\mathbf{x} \sim p(\mathbf{x}|\boldsymbol{\theta})$;
- $\mathbf{x} \sim \pi(\mathbf{x})$, $\mathbf{z} = \mathbf{f}_{\theta}(\mathbf{x})$, so $\mathbf{z} \sim p(\mathbf{z}|\boldsymbol{\theta})$;

*Papamakarios G., Pavlakou T., Murray I. Masked Autoregressive Flow for Density Estimation, 2017*

# Normalizing Flows KL Duality

### Theorem

$$\arg\min_{\boldsymbol{\theta}} \mathrm{KL}(\pi(\mathbf{x})\|p(\mathbf{x}|\boldsymbol{\theta})) = \arg\min_{\boldsymbol{\theta}} \mathrm{KL}(p(\mathbf{z}|\boldsymbol{\theta})\|p(\mathbf{z})).$$

### Proof

- $\mathbf{z} \sim p(\mathbf{z})$, $\mathbf{x} = \mathbf{g}_{\boldsymbol{\theta}}(\mathbf{z})$, so $\mathbf{x} \sim p(\mathbf{x}|\boldsymbol{\theta})$;
- $\mathbf{x} \sim \pi(\mathbf{x})$, $\mathbf{z} = \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})$, so $\mathbf{z} \sim p(\mathbf{z}|\boldsymbol{\theta})$;

$$\log p(\mathbf{z}|\boldsymbol{\theta}) = \log \pi(\mathbf{g}_{\boldsymbol{\theta}}(\mathbf{z})) + \log|\det(\mathbf{J_g})|;$$

$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \log p(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})) + \log|\det(\mathbf{J_f})|.$$

*Papamakarios G., Pavlakou T., Murray I. Masked Autoregressive Flow for Density Estimation, 2017*

# Normalizing Flows KL Duality

## Theorem

$$\arg\min_{\boldsymbol{\theta}} \mathrm{KL}(\pi(\mathbf{x})\|p(\mathbf{x}|\boldsymbol{\theta})) = \arg\min_{\boldsymbol{\theta}} \mathrm{KL}(p(\mathbf{z}|\boldsymbol{\theta})\|p(\mathbf{z})).$$

## Proof

- $\mathbf{z} \sim p(\mathbf{z})$, $\mathbf{x} = \mathbf{g}_{\boldsymbol{\theta}}(\mathbf{z})$, so $\mathbf{x} \sim p(\mathbf{x}|\boldsymbol{\theta})$;
- $\mathbf{x} \sim \pi(\mathbf{x})$, $\mathbf{z} = \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})$, so $\mathbf{z} \sim p(\mathbf{z}|\boldsymbol{\theta})$;

$$\log p(\mathbf{z}|\boldsymbol{\theta}) = \log \pi(\mathbf{g}_{\boldsymbol{\theta}}(\mathbf{z})) + \log|\det(\mathbf{J_g})|;$$
$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \log p(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})) + \log|\det(\mathbf{J_f})|.$$

$$\mathrm{KL}\left(p(\mathbf{z}|\boldsymbol{\theta})\|p(\mathbf{z})\right) = \mathbb{E}_{p(\mathbf{z}|\boldsymbol{\theta})}\left[\log p(\mathbf{z}|\boldsymbol{\theta}) - \log p(\mathbf{z})\right]$$

Papamakarios G., Pavlakou T., Murray I. *Masked Autoregressive Flow for Density Estimation*, 2017

# Normalizing Flows KL Duality

## Theorem

$$\arg\min_{\boldsymbol{\theta}} \mathrm{KL}(\pi(\mathbf{x})\|p(\mathbf{x}|\boldsymbol{\theta})) = \arg\min_{\boldsymbol{\theta}} \mathrm{KL}(p(\mathbf{z}|\boldsymbol{\theta})\|p(\mathbf{z})).$$

## Proof

▶ $\mathbf{z} \sim p(\mathbf{z})$, $\mathbf{x} = \mathbf{g}_{\boldsymbol{\theta}}(\mathbf{z})$, so $\mathbf{x} \sim p(\mathbf{x}|\boldsymbol{\theta})$;

▶ $\mathbf{x} \sim \pi(\mathbf{x})$, $\mathbf{z} = \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})$, so $\mathbf{z} \sim p(\mathbf{z}|\boldsymbol{\theta})$;

$$\log p(\mathbf{z}|\boldsymbol{\theta}) = \log \pi(\mathbf{g}_{\boldsymbol{\theta}}(\mathbf{z})) + \log|\det(\mathbf{J_g})|;$$
$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \log p(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})) + \log|\det(\mathbf{J_f})|.$$

$$\mathrm{KL}\left(p(\mathbf{z}|\boldsymbol{\theta})\|p(\mathbf{z})\right) = \mathbb{E}_{p(\mathbf{z}|\boldsymbol{\theta})}\left[\log p(\mathbf{z}|\boldsymbol{\theta}) - \log p(\mathbf{z})\right] =$$
$$= \mathbb{E}_{p(\mathbf{z}|\boldsymbol{\theta})}\left[\log \pi(\mathbf{g}_{\boldsymbol{\theta}}(\mathbf{z})) + \log|\det(\mathbf{J_g})| - \log p(\mathbf{z})\right]$$

*Papamakarios G., Pavlakou T., Murray I. Masked Autoregressive Flow for Density Estimation, 2017*

# Normalizing Flows KL Duality

## Theorem

$$\arg\min_{\boldsymbol{\theta}} \mathrm{KL}(\pi(\mathbf{x})\|p(\mathbf{x}|\boldsymbol{\theta})) = \arg\min_{\boldsymbol{\theta}} \mathrm{KL}(p(\mathbf{z}|\boldsymbol{\theta})\|p(\mathbf{z})).$$

## Proof

- $\mathbf{z} \sim p(\mathbf{z})$, $\mathbf{x} = \mathbf{g}_{\boldsymbol{\theta}}(\mathbf{z})$, so $\mathbf{x} \sim p(\mathbf{x}|\boldsymbol{\theta})$;
- $\mathbf{x} \sim \pi(\mathbf{x})$, $\mathbf{z} = \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})$, so $\mathbf{z} \sim p(\mathbf{z}|\boldsymbol{\theta})$;

$$\log p(\mathbf{z}|\boldsymbol{\theta}) = \log \pi(\mathbf{g}_{\boldsymbol{\theta}}(\mathbf{z})) + \log |\det(\mathbf{J_g})|;$$
$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \log p(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})) + \log |\det(\mathbf{J_f})|.$$

$$\mathrm{KL}\left(p(\mathbf{z}|\boldsymbol{\theta})\|p(\mathbf{z})\right) = \mathbb{E}_{p(\mathbf{z}|\boldsymbol{\theta})}\left[\log p(\mathbf{z}|\boldsymbol{\theta}) - \log p(\mathbf{z})\right] =$$
$$= \mathbb{E}_{p(\mathbf{z}|\boldsymbol{\theta})}\left[\log \pi(\mathbf{g}_{\boldsymbol{\theta}}(\mathbf{z})) + \log |\det(\mathbf{J_g})| - \log p(\mathbf{z})\right] =$$
$$= \mathbb{E}_{\pi(\mathbf{x})}\left[\log \pi(\mathbf{x}) - \log |\det(\mathbf{J_f})| - \log p(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}))\right]$$

Papamakarios G., Pavlakou T., Murray I. *Masked Autoregressive Flow for Density Estimation*, 2017

# Normalizing Flows KL Duality

### Theorem

$$\arg\min_{\boldsymbol{\theta}} \mathrm{KL}(\pi(\mathbf{x}) \| p(\mathbf{x}|\boldsymbol{\theta})) = \arg\min_{\boldsymbol{\theta}} \mathrm{KL}(p(\mathbf{z}|\boldsymbol{\theta}) \| p(\mathbf{z})).$$

### Proof

- $\mathbf{z} \sim p(\mathbf{z})$, $\mathbf{x} = \mathbf{g}_{\boldsymbol{\theta}}(\mathbf{z})$, so $\mathbf{x} \sim p(\mathbf{x}|\boldsymbol{\theta})$;
- $\mathbf{x} \sim \pi(\mathbf{x})$, $\mathbf{z} = \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})$, so $\mathbf{z} \sim p(\mathbf{z}|\boldsymbol{\theta})$;

$$\log p(\mathbf{z}|\boldsymbol{\theta}) = \log \pi(\mathbf{g}_{\boldsymbol{\theta}}(\mathbf{z})) + \log |\det(\mathbf{J_g})|;$$
$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \log p(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})) + \log |\det(\mathbf{J_f})|.$$

$$\mathrm{KL}\left(p(\mathbf{z}|\boldsymbol{\theta}) \| p(\mathbf{z})\right) = \mathbb{E}_{p(\mathbf{z}|\boldsymbol{\theta})}\left[\log p(\mathbf{z}|\boldsymbol{\theta}) - \log p(\mathbf{z})\right] =$$
$$= \mathbb{E}_{p(\mathbf{z}|\boldsymbol{\theta})}\left[\log \pi(\mathbf{g}_{\boldsymbol{\theta}}(\mathbf{z})) + \log |\det(\mathbf{J_g})| - \log p(\mathbf{z})\right] =$$
$$= \mathbb{E}_{\pi(\mathbf{x})}\left[\log \pi(\mathbf{x}) - \log |\det(\mathbf{J_f})| - \log p(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}))\right] =$$
$$= \mathbb{E}_{\pi(\mathbf{x})}\left[\log \pi(\mathbf{x}) - \log p(\mathbf{x}|\boldsymbol{\theta})\right] = \mathrm{KL}(\pi(\mathbf{x}) \| p(\mathbf{x}|\boldsymbol{\theta})).$$

Papamakarios G., Pavlakou T., Murray I. *Masked Autoregressive Flow for Density Estimation*, 2017

# Outline

# Bayesian Framework

## Bayes' Theorem

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

- ▶ $\mathbf{x}$: observed variables;
- ▶ $\boldsymbol{\theta}$: unknown latent variables/parameters;
- ▶ $p(\mathbf{x}|\boldsymbol{\theta})$: likelihood;
- ▶ $p(\mathbf{x}) = \int p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta}$: evidence;
- ▶ $p(\boldsymbol{\theta})$: prior distribution;
- ▶ $p(\boldsymbol{\theta}|\mathbf{x})$: posterior distribution.

# Bayesian Framework

## Bayes' Theorem

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

- ▶ $\mathbf{x}$: observed variables;
- ▶ $\boldsymbol{\theta}$: unknown latent variables/parameters;
- ▶ $p(\mathbf{x}|\boldsymbol{\theta})$: likelihood;
- ▶ $p(\mathbf{x}) = \int p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta}$: evidence;
- ▶ $p(\boldsymbol{\theta})$: prior distribution;
- ▶ $p(\boldsymbol{\theta}|\mathbf{x})$: posterior distribution.

## Interpretation

- ▶ We begin with unknown variables $\boldsymbol{\theta}$ and a prior belief $p(\boldsymbol{\theta})$.
- ▶ Once data $\mathbf{x}$ is observed, the posterior $p(\boldsymbol{\theta}|\mathbf{x})$ incorporates both prior beliefs and evidence from the data.

# Bayesian Framework

Consider the case where the unobserved variables $\boldsymbol{\theta}$ are model parameters (i.e., $\boldsymbol{\theta}$ are random variables).

- $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^{n}$: observed samples;
- $p(\boldsymbol{\theta})$: prior distribution.

# Bayesian Framework

Consider the case where the unobserved variables $\boldsymbol{\theta}$ are model parameters (i.e., $\boldsymbol{\theta}$ are random variables).

- $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^{n}$: observed samples;
- $p(\boldsymbol{\theta})$: prior distribution.

## Posterior Distribution

$$p(\boldsymbol{\theta}|\mathbf{X}) = \frac{p(\mathbf{X}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{X})} = \frac{p(\mathbf{X}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{\int p(\mathbf{X}|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta}}$$

# Bayesian Framework

Consider the case where the unobserved variables $\boldsymbol{\theta}$ are model parameters (i.e., $\boldsymbol{\theta}$ are random variables).

- $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^{n}$: observed samples;
- $p(\boldsymbol{\theta})$: prior distribution.

## Posterior Distribution

$$p(\boldsymbol{\theta}|\mathbf{X}) = \frac{p(\mathbf{X}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{X})} = \frac{p(\mathbf{X}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{\int p(\mathbf{X}|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta}}$$

If the evidence $p(\mathbf{X})$ is intractable (due to high-dimensional integration), the posterior cannot be computed exactly.

# Bayesian Framework

Consider the case where the unobserved variables $\boldsymbol{\theta}$ are model parameters (i.e., $\boldsymbol{\theta}$ are random variables).

- $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n$: observed samples;
- $p(\boldsymbol{\theta})$: prior distribution.

## Posterior Distribution

$$p(\boldsymbol{\theta}|\mathbf{X}) = \frac{p(\mathbf{X}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{X})} = \frac{p(\mathbf{X}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{\int p(\mathbf{X}|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta}}$$

If the evidence $p(\mathbf{X})$ is intractable (due to high-dimensional integration), the posterior cannot be computed exactly.

## Maximum a Posteriori (MAP) Estimation

$$\boldsymbol{\theta}^* = \arg\max_{\boldsymbol{\theta}} p(\boldsymbol{\theta}|\mathbf{X}) = \arg\max_{\boldsymbol{\theta}}(\log p(\mathbf{X}|\boldsymbol{\theta}) + \log p(\boldsymbol{\theta}))$$

# Latent Variable Models (LVM)

## Maximum Likelihood Extimation (MLE) Problem

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

# Latent Variable Models (LVM)

### Maximum Likelihood Extimation (MLE) Problem

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

The distribution $p(\mathbf{x}|\boldsymbol{\theta})$ can be highly complex and often intractable (just like the true data distribution $\pi(\mathbf{x})$).

# Latent Variable Models (LVM)

## Maximum Likelihood Extimation (MLE) Problem

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

The distribution $p(\mathbf{x}|\boldsymbol{\theta})$ can be highly complex and often intractable (just like the true data distribution $\pi(\mathbf{x})$).

## Extended Probabilistic Model
Introduce a latent variable $\mathbf{z}$ for each observed sample $\mathbf{x}$:

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

# Latent Variable Models (LVM)

### Maximum Likelihood Extimation (MLE) Problem

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

The distribution $p(\mathbf{x}|\boldsymbol{\theta})$ can be highly complex and often intractable (just like the true data distribution $\pi(\mathbf{x})$).

### Extended Probabilistic Model

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

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

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

# Latent Variable Models (LVM)

### Maximum Likelihood Extimation (MLE) Problem

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

The distribution $p(\mathbf{x}|\boldsymbol{\theta})$ can be highly complex and often intractable (just like the true data distribution $\pi(\mathbf{x})$).

### Extended Probabilistic Model

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

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

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

### Motivation

Both $p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta})$ and $p(\mathbf{z})$ are usually much simpler than $p(\mathbf{x}|\boldsymbol{\theta})$.

# Latent Variable Models (LVM)

$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \log \int p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta})p(\mathbf{z})d\mathbf{z} \to \max_{\boldsymbol{\theta}}$$

Bishop C. Pattern Recognition and Machine Learning, 2006

# Latent Variable Models (LVM)

$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \log \int p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta}) p(\mathbf{z}) d\mathbf{z} \rightarrow \max_{\boldsymbol{\theta}}$$
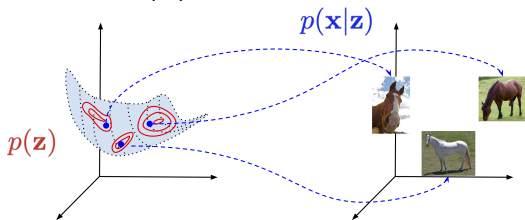
Examples

*Mixture of Gaussians*
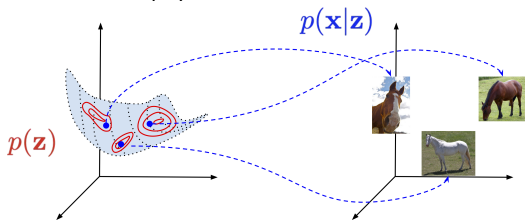


- ▶ $p(\mathbf{x}|z, \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z)$
- ▶ $p(z) = \mathrm{Categorical}(\boldsymbol{\pi})$

---

*Bishop C. Pattern Recognition and Machine Learning, 2006*

# Latent Variable Models (LVM)

$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \log \int p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta}) p(\mathbf{z}) d\mathbf{z} \to \max_{\boldsymbol{\theta}}$$

## Examples

*Mixture of Gaussians*                    *PCA Model*



- $p(\mathbf{x}|z, \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z)$
- $p(z) = \mathrm{Categorical}(\boldsymbol{\pi})$

- $p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{Wz} + \boldsymbol{\mu}, \sigma^2 \mathbf{I})$
- $p(\mathbf{z}) = \mathcal{N}(0, \mathbf{I})$

---

*Bishop C. Pattern Recognition and Machine Learning, 2006*

# MLE for LVM

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

# MLE for LVM

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

# MLE for LVM

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



## Naive Approach

$$p(\mathbf{x}|\boldsymbol{\theta}) = \int p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta})p(\mathbf{z})d\mathbf{z} = \mathbb{E}_{p(\mathbf{z})}p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta}) \approx \frac{1}{K}\sum_{k=1}^{K} p(\mathbf{x}|\mathbf{z}_k, \boldsymbol{\theta}),$$

where $\mathbf{z}_k \sim p(\mathbf{z})$.

---

*image credit: https://jmtomczak.github.io/blog/4/4_VAE.html*

# MLE for LVM

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



## Naive Approach

$$p(\mathbf{x} | \boldsymbol{\theta}) = \int p(\mathbf{x} | \mathbf{z}, \boldsymbol{\theta}) p(\mathbf{z}) d\mathbf{z} = \mathbb{E}_{p(\mathbf{z})} p(\mathbf{x} | \mathbf{z}, \boldsymbol{\theta}) \approx \frac{1}{K} \sum_{k=1}^{K} p(\mathbf{x} | \mathbf{z}_k, \boldsymbol{\theta}),$$

where $\mathbf{z}_k \sim p(\mathbf{z})$.

**Challenge:** As the dimensionality of $\mathbf{z}$ increases, the number of samples needed to adequately cover the latent space grows exponentially.

# Outline

# ELBO Derivation I

Inequality Derivation

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

# ELBO Derivation I

### Inequality Derivation

$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \log \int p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}) d\mathbf{z} = \log \int \frac{q(\mathbf{z})}{q(\mathbf{z})} p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}) d\mathbf{z}$$

# ELBO Derivation I

### Inequality Derivation

$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \log \int p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}) d\mathbf{z} = \log \int \frac{q(\mathbf{z})}{q(\mathbf{z})} p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}) d\mathbf{z} =$$
$$= \log \mathbb{E}_q \left[ \frac{p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta})}{q(\mathbf{z})} \right]$$

# ELBO Derivation I

### Inequality Derivation

$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \log \int p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta})d\mathbf{z} = \log \int \frac{q(\mathbf{z})}{q(\mathbf{z})}p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta})d\mathbf{z} =$$

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

# ELBO Derivation I

### Inequality Derivation

$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \log \int p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}) d\mathbf{z} = \log \int \frac{q(\mathbf{z})}{q(\mathbf{z})} p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}) d\mathbf{z} =$$

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

Here, $q(\mathbf{z})$ is any distribution such that $\int q(\mathbf{z}) d\mathbf{z} = 1$.

# ELBO Derivation I

### Inequality Derivation

$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \log \int p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}) d\mathbf{z} = \log \int \frac{q(\mathbf{z})}{q(\mathbf{z})} p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}) d\mathbf{z} =$$
$$= \log \mathbb{E}_q \left[ \frac{p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta})}{q(\mathbf{z})} \right] \geq \mathbb{E}_q \log \frac{p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta})}{q(\mathbf{z})} = \mathcal{L}_{q,\theta}(\mathbf{x})$$

Here, $q(\mathbf{z})$ is any distribution such that $\int q(\mathbf{z}) d\mathbf{z} = 1$.

### Variational Evidence Lower Bound (ELBO)

$$\mathcal{L}_{q,\theta}(\mathbf{x}) = \mathbb{E}_q \log \frac{p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta})}{q(\mathbf{z})} \leq \log p(\mathbf{x}|\boldsymbol{\theta})$$

# ELBO Derivation I

### Inequality Derivation

$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \log \int p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta})d\mathbf{z} = \log \int \frac{q(\mathbf{z})}{q(\mathbf{z})} p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta})d\mathbf{z} =$$

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

Here, $q(\mathbf{z})$ is any distribution such that $\int q(\mathbf{z})d\mathbf{z} = 1$.

### Variational Evidence Lower Bound (ELBO)

$$\mathcal{L}_{q,\theta}(\mathbf{x}) = \mathbb{E}_q \log \frac{p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta})}{q(\mathbf{z})} \leq \log p(\mathbf{x}|\boldsymbol{\theta})$$

This inequality holds for any choice of $q(\mathbf{z})$.

# ELBO Derivation II

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

## Equality Derivation

$$\mathcal{L}_{q,\boldsymbol{\theta}}(\mathbf{x}) = \int q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta})}{q(\mathbf{z})} d\mathbf{z}$$

# ELBO Derivation II

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

Equality Derivation

$$\mathcal{L}_{q,\boldsymbol{\theta}}(\mathbf{x}) = \int q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta})}{q(\mathbf{z})} d\mathbf{z} =$$
$$= \int q(\mathbf{z}) \log \frac{p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) p(\mathbf{x}|\boldsymbol{\theta})}{q(\mathbf{z})} d\mathbf{z}$$

# ELBO Derivation II

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

## Equality Derivation

$$\mathcal{L}_{q,\boldsymbol{\theta}}(\mathbf{x}) = \int q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta})}{q(\mathbf{z})} d\mathbf{z} =$$

$$= \int q(\mathbf{z}) \log \frac{p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})p(\mathbf{x}|\boldsymbol{\theta})}{q(\mathbf{z})} d\mathbf{z} =$$

$$= \int q(\mathbf{z}) \log p(\mathbf{x}|\boldsymbol{\theta}) d\mathbf{z} + \int q(\mathbf{z}) \log \frac{p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})}{q(\mathbf{z})} d\mathbf{z}$$

# ELBO Derivation II

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

Equality Derivation

$$\begin{aligned}
\mathcal{L}_{q,\boldsymbol{\theta}}(\mathbf{x}) &= \int q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta})}{q(\mathbf{z})} d\mathbf{z} = \\
&= \int q(\mathbf{z}) \log \frac{p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})p(\mathbf{x}|\boldsymbol{\theta})}{q(\mathbf{z})} d\mathbf{z} = \\
&= \int q(\mathbf{z}) \log p(\mathbf{x}|\boldsymbol{\theta}) d\mathbf{z} + \int q(\mathbf{z}) \log \frac{p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})}{q(\mathbf{z})} d\mathbf{z} = \\
&= \log p(\mathbf{x}|\boldsymbol{\theta}) - \mathrm{KL}(q(\mathbf{z})\|p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}))
\end{aligned}$$

# ELBO Derivation II

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

Equality Derivation

$$\mathcal{L}_{q,\boldsymbol{\theta}}(\mathbf{x}) = \int q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta})}{q(\mathbf{z})} d\mathbf{z} =$$

$$= \int q(\mathbf{z}) \log \frac{p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) p(\mathbf{x}|\boldsymbol{\theta})}{q(\mathbf{z})} d\mathbf{z} =$$

$$= \int q(\mathbf{z}) \log p(\mathbf{x}|\boldsymbol{\theta}) d\mathbf{z} + \int q(\mathbf{z}) \log \frac{p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})}{q(\mathbf{z})} d\mathbf{z} =$$

$$= \log p(\mathbf{x}|\boldsymbol{\theta}) - \mathrm{KL}(q(\mathbf{z})\|p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}))$$

Variational Decomposition

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

# ELBO Derivation II

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

Equality Derivation

$$\mathcal{L}_{q,\boldsymbol{\theta}}(\mathbf{x}) = \int q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta})}{q(\mathbf{z})} d\mathbf{z} =$$

$$= \int q(\mathbf{z}) \log \frac{p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) p(\mathbf{x}|\boldsymbol{\theta})}{q(\mathbf{z})} d\mathbf{z} =$$

$$= \int q(\mathbf{z}) \log p(\mathbf{x}|\boldsymbol{\theta}) d\mathbf{z} + \int q(\mathbf{z}) \log \frac{p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})}{q(\mathbf{z})} d\mathbf{z} =$$

$$= \log p(\mathbf{x}|\boldsymbol{\theta}) - \mathrm{KL}(q(\mathbf{z})\|p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}))$$

Variational Decomposition

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

Here, $\mathrm{KL}(q(\mathbf{z})\|p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})) \geq 0$.

# Variational Evidence Lower Bound (ELBO)

$$\mathcal{L}_{q,\theta}(\mathbf{x}) = \int q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z})} d\mathbf{z}$$

# Variational Evidence Lower Bound (ELBO)

$$\mathcal{L}_{q,\theta}(\mathbf{x}) = \int q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z})} d\mathbf{z}$$

$$= \int q(\mathbf{z}) \log p(\mathbf{x}|\mathbf{z}, \theta) d\mathbf{z} + \int q(\mathbf{z}) \log \frac{p(\mathbf{z})}{q(\mathbf{z})} d\mathbf{z}$$

# Variational Evidence Lower Bound (ELBO)

$$\mathcal{L}_{q,\theta}(\mathbf{x}) = \int q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z})} d\mathbf{z}$$

$$= \int q(\mathbf{z}) \log p(\mathbf{x}|\mathbf{z}, \theta) d\mathbf{z} + \int q(\mathbf{z}) \log \frac{p(\mathbf{z})}{q(\mathbf{z})} d\mathbf{z}$$

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

# Variational Evidence Lower Bound (ELBO)

$$\mathcal{L}_{q,\theta}(\mathbf{x}) = \int q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z} | \theta)}{q(\mathbf{z})} d\mathbf{z}$$

$$= \int q(\mathbf{z}) \log p(\mathbf{x} | \mathbf{z}, \theta) d\mathbf{z} + \int q(\mathbf{z}) \log \frac{p(\mathbf{z})}{q(\mathbf{z})} d\mathbf{z}$$

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

Log-Likelihood Decomposition

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

# Variational Evidence Lower Bound (ELBO)

$$\mathcal{L}_{q,\theta}(\mathbf{x}) = \int q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta})}{q(\mathbf{z})} d\mathbf{z}$$

$$= \int q(\mathbf{z}) \log p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta}) d\mathbf{z} + \int q(\mathbf{z}) \log \frac{p(\mathbf{z})}{q(\mathbf{z})} d\mathbf{z}$$

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

Log-Likelihood Decomposition

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

$$= \mathbb{E}_q \log p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta}) - \mathrm{KL}(q(\mathbf{z})\|p(\mathbf{z})) + \mathrm{KL}(q(\mathbf{z})\|p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})).$$

# Variational Evidence Lower Bound (ELBO)

$$\mathcal{L}_{q,\theta}(\mathbf{x}) = \int q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z})} d\mathbf{z}$$

$$= \int q(\mathbf{z}) \log p(\mathbf{x}|\mathbf{z}, \theta) d\mathbf{z} + \int q(\mathbf{z}) \log \frac{p(\mathbf{z})}{q(\mathbf{z})} d\mathbf{z}$$

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

## Log-Likelihood Decomposition

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

▶ Instead of maximizing the likelihood, maximize the ELBO:

$$\max_{\theta} p(\mathbf{x}|\theta) \quad \to \quad \max_{q,\theta} \mathcal{L}_{q,\theta}(\mathbf{x})$$

# Variational Evidence Lower Bound (ELBO)

$$\mathcal{L}_{q,\theta}(\mathbf{x}) = \int q(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z} | \theta)}{q(\mathbf{z})} d\mathbf{z}$$

$$= \int q(\mathbf{z}) \log p(\mathbf{x} | \mathbf{z}, \theta) d\mathbf{z} + \int q(\mathbf{z}) \log \frac{p(\mathbf{z})}{q(\mathbf{z})} d\mathbf{z}$$

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

## Log-Likelihood Decomposition

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

▶ Instead of maximizing the likelihood, maximize the ELBO:

$$\max_{\theta} p(\mathbf{x} | \theta) \quad \rightarrow \quad \max_{q, \theta} \mathcal{L}_{q,\theta}(\mathbf{x})$$

▶ Maximizing the ELBO with respect to the **variational** distribution $q$ is equivalent to minimizing the KL divergence:

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

# Outline

# EM-Algorithm

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

$$= \mathbb{E}_q \Big[ \log p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta}) - \log \frac{q(\mathbf{z})}{p(\mathbf{z})} \Big] d\mathbf{z} \to \max_{q,\boldsymbol{\theta}}.$$

# EM-Algorithm

$$\mathcal{L}_{q,\theta}(\mathbf{x}) = \mathbb{E}_q \log p(\mathbf{x}|\mathbf{z}, \theta) - \mathrm{KL}(q(\mathbf{z})\|p(\mathbf{z})) =$$
$$= \mathbb{E}_q \Big[ \log p(\mathbf{x}|\mathbf{z}, \theta) - \log \frac{q(\mathbf{z})}{p(\mathbf{z})} \Big] d\mathbf{z} \to \max_{q, \theta}.$$

Block-Coordinate Optimization

▶ Initialize $\theta^*$;

# EM-Algorithm

$$\mathcal{L}_{q,\boldsymbol{\theta}}(\mathbf{x}) = \mathbb{E}_q \log p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta}) - \mathrm{KL}(q(\mathbf{z})\|p(\mathbf{z})) =$$
$$= \mathbb{E}_q \left[ \log p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta}) - \log \frac{q(\mathbf{z})}{p(\mathbf{z})} \right] d\mathbf{z} \to \max_{q,\boldsymbol{\theta}}.$$

Block-Coordinate Optimization

- Initialize $\boldsymbol{\theta}^*$;
- **E-step** (optimize $\mathcal{L}_{q,\boldsymbol{\theta}}(\mathbf{x})$ over $q$):
  $$q^*(\mathbf{z}) = \arg\max_q \mathcal{L}_{q,\boldsymbol{\theta}^*}(\mathbf{x}) =$$
  $$= \arg\min_q \mathrm{KL}(q(\mathbf{z})\|p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}^*)) = p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}^*);$$

# EM-Algorithm

$$\mathcal{L}_{q,\boldsymbol{\theta}}(\mathbf{x}) = \mathbb{E}_q \log p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta}) - \mathrm{KL}(q(\mathbf{z})\|p(\mathbf{z})) =$$
$$= \mathbb{E}_q \Big[ \log p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta}) - \log \frac{q(\mathbf{z})}{p(\mathbf{z})} \Big] d\mathbf{z} \to \max_{q,\boldsymbol{\theta}}.$$

## Block-Coordinate Optimization

▶ Initialize $\boldsymbol{\theta}^*$;

▶ **E-step** (optimize $\mathcal{L}_{q,\boldsymbol{\theta}}(\mathbf{x})$ over $q$):
$$q^*(\mathbf{z}) = \arg\max_q \mathcal{L}_{q,\boldsymbol{\theta}^*}(\mathbf{x}) =$$
$$= \arg\min_q \mathrm{KL}(q(\mathbf{z})\|p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}^*)) = p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}^*);$$

▶ **M-step** (optimize $\mathcal{L}_{q,\boldsymbol{\theta}}(\mathbf{x})$ over $\boldsymbol{\theta}$):
$$\boldsymbol{\theta}^* = \arg\max_{\boldsymbol{\theta}} \mathcal{L}_{q^*,\boldsymbol{\theta}}(\mathbf{x});$$

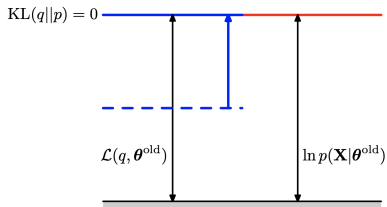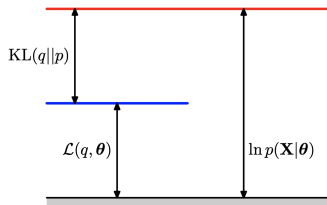# EM-Algorithm

$$\mathcal{L}_{q,\theta}(\mathbf{x}) = \mathbb{E}_q \log p(\mathbf{x}|\mathbf{z}, \theta) - \mathrm{KL}(q(\mathbf{z})\|p(\mathbf{z})) =$$
$$= \mathbb{E}_q \Big[ \log p(\mathbf{x}|\mathbf{z}, \theta) - \log \frac{q(\mathbf{z})}{p(\mathbf{z})} \Big] d\mathbf{z} \to \max_{q,\theta}.$$
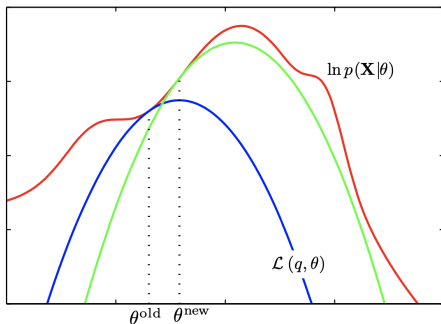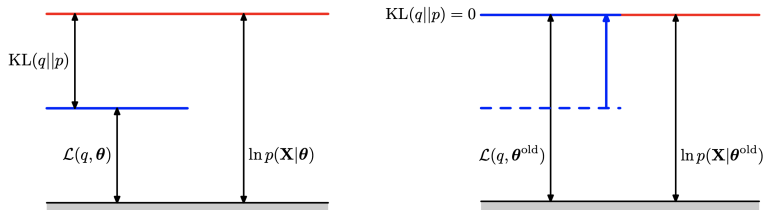
## Block-Coordinate Optimization

▶ Initialize $\theta^*$;

▶ **E-step** (optimize $\mathcal{L}_{q,\theta}(\mathbf{x})$ over $q$):
$$q^*(\mathbf{z}) = \arg\max_q \mathcal{L}_{q,\theta^*}(\mathbf{x}) =$$
$$= \arg\min_q \mathrm{KL}(q(\mathbf{z})\|p(\mathbf{z}|\mathbf{x}, \theta^*)) = p(\mathbf{z}|\mathbf{x}, \theta^*);$$

▶ **M-step** (optimize $\mathcal{L}_{q,\theta}(\mathbf{x})$ over $\theta$):
$$\theta^* = \arg\max_\theta \mathcal{L}_{q^*,\theta}(\mathbf{x});$$

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

# EM-Algorithm Illustration



Bishop C. Pattern Recognition and Machine Learning, 2006

# EM-Algorithm Illustration



Bishop C. Pattern Recognition and Machine Learning, 2006

# Summary

- ▶ Flow duality establishes the relationship between the data and latent spaces using forward and reverse KL formulations.

- ▶ The Bayesian framework generalizes nearly all standard machine learning methods.

- ▶ LVMs introduce latent representations for observed data, enabling more interpretable models.

- ▶ LVMs maximize the variational evidence lower bound (ELBO) to obtain maximum likelihood estimates for the parameters.

- ▶ The general variational EM algorithm optimizes the ELBO within LVMs to recover the MLE for the parameters $\boldsymbol{\theta}$.