

# Deep Generative Models

## Lecture 2

Roman Isachenko



Ozon Masters

Spring, 2021

## Recap of previous lecture

We are given i.i.d. samples  $\{\mathbf{x}_i\}_{i=1}^n \in X$  (e.g.  $X = \mathbb{R}^m$ ) from unknown distribution  $\pi(\mathbf{x})$ .

### Goal

We would like to learn a distribution  $\pi(\mathbf{x})$  for

- ▶ evaluating  $\pi(\mathbf{x})$  for new samples (how likely to get object  $\mathbf{x}$ ?);
- ▶ sampling from  $\pi(\mathbf{x})$  (to get new objects  $\mathbf{x} \sim \pi(\mathbf{x})$ ).

### Challenge

Data is complex and high-dimensional.

### MLE problem

Fix probabilistic model  $p(\mathbf{x}|\boldsymbol{\theta})$  – the set of parameterized distributions, such that  $p(\mathbf{x}|\boldsymbol{\theta}) \approx \pi(\mathbf{x})$ .

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

# Recap of previous lecture

## Likelihood as product of conditionals

Let  $\mathbf{x} = (x_1, \dots, x_m)$ ,  $\mathbf{x}_{1:i} = (x_1, \dots, x_i)$ . Then

$$p(\mathbf{x}|\boldsymbol{\theta}) = \prod_{i=1}^m p(x_i|\mathbf{x}_{1:i-1}, \boldsymbol{\theta}); \quad \log p(\mathbf{x}|\boldsymbol{\theta}) = \sum_{i=1}^m \log p(x_i|\mathbf{x}_{1:i-1}, \boldsymbol{\theta}).$$

## MLE problem for autoregressive model

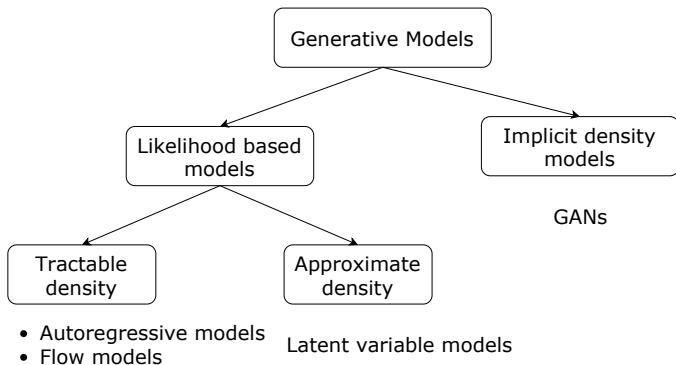
$$\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} p(\mathbf{X}|\boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} \prod_{i=1}^n p(\mathbf{x}_i|\boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^n \log p(\mathbf{x}_i|\boldsymbol{\theta}).$$

## Sampling

$$\hat{x}_1 \sim p(x_1|\boldsymbol{\theta}), \quad \hat{x}_2 \sim p(x_2|\hat{x}_1, \boldsymbol{\theta}), \dots, \quad \hat{x}_m \sim p(x_m|\hat{\mathbf{x}}_{1:m-1}, \boldsymbol{\theta})$$

New generated object is  $\hat{\mathbf{x}} = (\hat{x}_1, \hat{x}_2, \dots, \hat{x}_m)$ .

# Generative models zoo



# 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;
- ▶  $p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{t})p(\mathbf{t})d\mathbf{t}$  – evidence;
- ▶  $p(\mathbf{t})$  – prior distribution,  $p(\mathbf{t}|\mathbf{x})$  – posterior distribution.

## Meaning

We have the prior knowledge  $p(\mathbf{t})$  of our unobserved variables  $\mathbf{t}$ . We have got the observed data  $\mathbf{x}$ . Posterior distribution  $p(\mathbf{t}|\mathbf{x})$  summarizes what you know after the data has been observed.

## Bayesian framework

Let consider the case, where the unobserved variables  $\mathbf{t}$  is our model parameters  $\theta$ .

- ▶  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n$  – observed samples;
- ▶  $p(\theta)$  – prior parameters distribution (we treat model parameters  $\theta$  as random variable).

## Posterior distribution

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

## Bayesian inference

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

Note the difference from

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

# Bayesian framework

## Posterior distribution

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

## Bayesian inference

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

If evidence  $p(\mathbf{X})$  is intractable (due to multidimensional integral), we can't get posterior distribution and make the honest inference.

## Maximum a posteriori (MAP) estimation

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

# Bayesian framework

## MAP estimation

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

Now  $\theta^*$  is a deterministic variable, but we could treat it as a random variable with density  $p(\theta|\mathbf{X}) = \delta(\theta - \theta^*)$ .

## Dirac delta function

$$\delta(x) = \begin{cases} +\infty, & x = 0; \\ 0, & x \neq 0; \end{cases} \quad \int f(x)\delta(x-y)dx = f(y).$$

## MAP inference

$$p(\mathbf{x}|\mathbf{X}) = \int p(\mathbf{x}|\theta)p(\theta|\mathbf{X})d\theta \approx p(\mathbf{x}|\theta^*).$$



# Latent variable models (LVM)

## 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.

## Extend probabilistic model

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

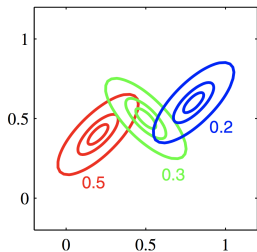
The distributions  $p(\mathbf{x}|\mathbf{z}, \theta)$  and  $p(\mathbf{z})$  could be quite simple.

# Latent variable models (LVM)

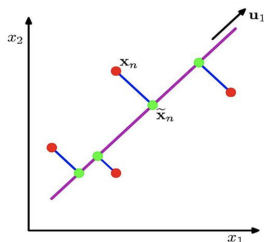
$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \log \int p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta})p(\mathbf{z})d\mathbf{z} \rightarrow \max_{\boldsymbol{\theta}}$$

## Examples

*Mixture of gaussians*



*PCA model*

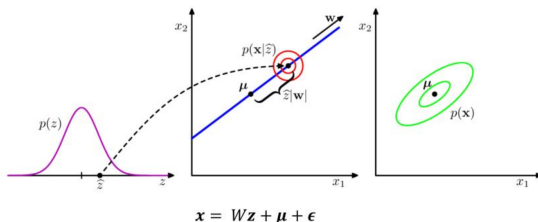


- ▶  $p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_{\mathbf{z}}, \boldsymbol{\Sigma}_{\mathbf{z}})$
- ▶  $p(\mathbf{z}) = \text{Categorical}(\mathbf{z}|\boldsymbol{\pi})$
- ▶  $p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \boldsymbol{\Sigma}_{\mathbf{z}})$
- ▶  $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$

# Latent variable models (LVM)

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

**PCA goal:** Project original data  $\mathbf{X}$  onto low latent space while maximizing the variance of the projected data.



- ▶  $p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \boldsymbol{\Sigma}_z)$
- ▶  $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$

# Incomplete likelihood

## MLE

$$\begin{aligned}\theta^* &= \arg \max_{\theta} p(\mathbf{X}, \mathbf{Z} | \theta) = \arg \max_{\theta} \prod_{i=1}^n p(\mathbf{x}_i, \mathbf{z}_i | \theta) = \\ &= \arg \max_{\theta} \sum_{i=1}^n \log p(\mathbf{x}_i, \mathbf{z}_i | \theta).\end{aligned}$$

Since  $\mathbf{Z}$  is unknown, maximize **incomplete likelihood**.

## MILE problem

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

## Variational lower bound

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

## Kullback-Leibler divergence

- ▶  $KL(q||p) = \int q(\mathbf{z}) \log \frac{q(\mathbf{z})}{p(\mathbf{z})} d\mathbf{z};$
- ▶  $KL(q||p) \geq 0;$
- ▶  $KL(q||p) = 0 \Leftrightarrow q \equiv p.$

## Variational lower bound

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

## Evidence Lower Bound (ELBO)

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

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

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

# EM-algorithm

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

## Block-coordinate optimization

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

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

- ▶ M-step

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

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

Ugly pic



# Amortized variational inference

## E-step

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

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

## Idea

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

## Variational Bayes

- ▶ E-step

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

- ▶ M-step

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

# 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

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

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

- ▶ M-step

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

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

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

**Difficulty:** number of samples  $n$ .

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

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

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, $\nabla_{\phi} \mathcal{L}(\phi, \theta)$ )

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

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:

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

First step is not an expectation due to the gradient.

### Possible solutions

- ▶ log-derivative trick;
- ▶ reparametrization trick.

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

Log-derivative trick

$$\nabla_{\xi} q(\eta|\xi) = q(\eta|\xi) \left( \frac{\nabla_{\xi} q(\eta|\xi)}{q(\eta|\xi)} \right) = q(\eta|\xi) \nabla_{\xi} \log q(\eta|\xi).$$

$$\nabla_{\phi} q(\mathbf{z}|\mathbf{x}, \phi) = q(\mathbf{z}|\mathbf{x}, \phi) \nabla_{\phi} \log q(\mathbf{z}|\mathbf{x}, \phi).$$

ELBO

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

Problem

Unstable solution with huge variance.

# 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  $r(\epsilon)$ .

$$\begin{aligned} \nabla_{\xi} \int q(\eta|\xi) h(\eta) d\eta &= \nabla_{\xi} \int r(\epsilon) h(g(\xi, \epsilon)) d\epsilon \\ &\approx \nabla_{\xi} h(g(\xi, \epsilon^*)), \quad \epsilon^* \sim r(\epsilon). \end{aligned}$$

## Example

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

## 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 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 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.

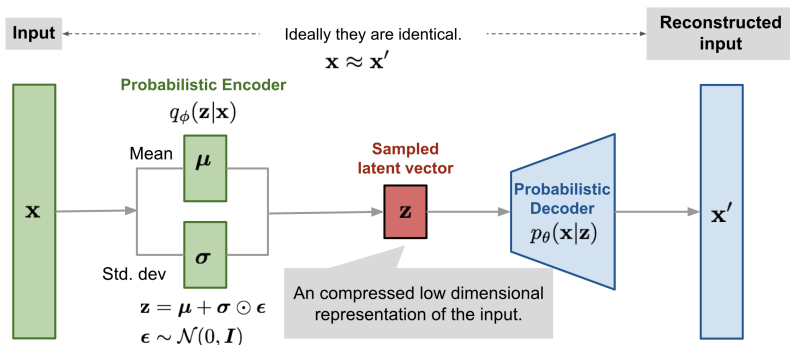
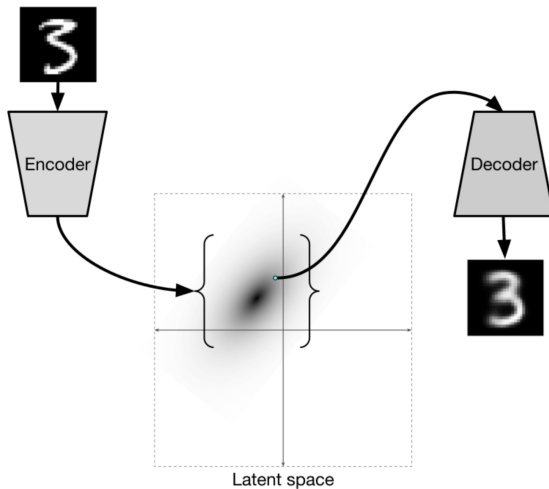


image credit:

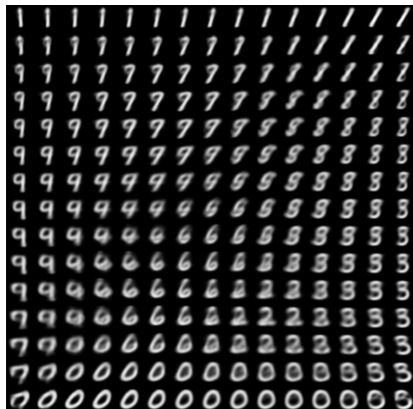
<https://lilianweng.github.io/lil-log/2018/08/12/from-autoencoder-to-beta-vae.html>

# Variational Autoencoder



# Variational Autoencoder

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



# 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.
- ▶ Latent variable models introduce latent representation of observed samples to make model more interpretable.
- ▶ To find MLE of model parameters LVM maximizes variational evidence lower bound.
- ▶ Maximization of ELBO is performed by general variational EM algorithm.
- ▶ Amortized inference allows to efficiently compute stochastic gradients for ELBO and use deep neural network for  $q(\mathbf{z}|\mathbf{x}, \phi)$  and  $p(\mathbf{x}|\mathbf{z}_i, \theta)$ .
- ▶ VAE model is a LVM model with encoder network for  $q(\mathbf{z}|\mathbf{x}, \phi)$  and decoder network for  $p(\mathbf{x}|\mathbf{z}_i, \theta)$ .