

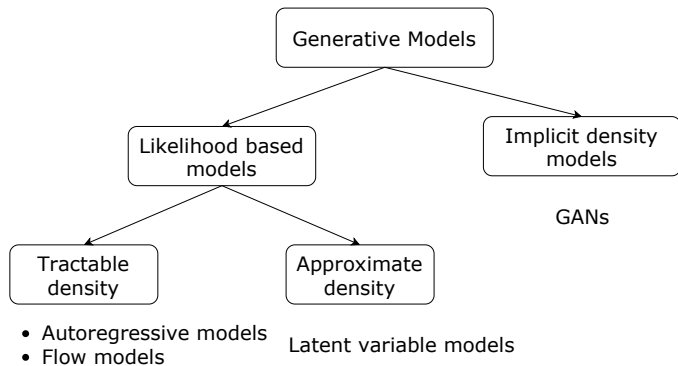
# Deep Generative Models

Roman Isachenko

Moscow Institute of Physics and Technology

2019

# Generative models zoo



# Bayesian framework

- ▶  $\mathbf{x}$  – observed samples;
- ▶  $\mathbf{z}$  – unobserved (latent) variables;
- ▶  $\theta$  – model parameters.

Discriminative

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

Classification/Regression

Generative

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

Generation of new samples  $(\mathbf{z}, \mathbf{x})$

# Bayesian framework

## Bayes theorem

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

## Full Bayesian inference

$$p(\mathbf{z}^*, \mathbf{x}^*|\mathbf{X}, \mathbf{Z}) = \int p(\mathbf{z}, \mathbf{x}|\theta)p(\theta|\mathbf{X}, \mathbf{Z})d\theta$$

## Maximum a posteriori (MAP)

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

# Latent variable models

## 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  $\mathbf{x}$

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

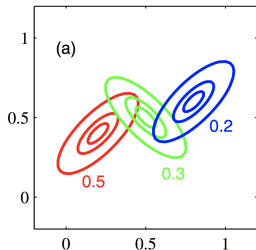
$$p(\mathbf{x}|\theta) = \int p(\mathbf{x}, \mathbf{z}|\theta) d\mathbf{z} = \int p(\mathbf{x}|\mathbf{z}, \theta)p(\mathbf{z}) d\mathbf{z}.$$

# Latent variable models

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

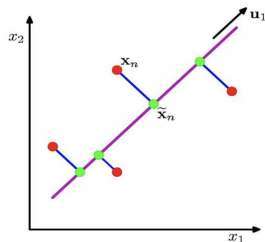
## Examples

*Mixture of gaussians*



- ▶  $p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_{\mathbf{z}}, \boldsymbol{\Sigma}_{\mathbf{z}}^2)$
- ▶  $p(\mathbf{z}) = \text{Cat}(\mathbf{z}|\boldsymbol{\pi})$

*PCA model*



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

# Incomplete likelihood

## MLE problem

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

Since  $Z$  is unknown, maximize **incomplete likelihood**.

## MILE problem

$$\begin{aligned}\theta^* &= \arg \max_{\theta} \log p(\mathbf{X}|\theta) = \arg \max_{\theta} \log \int p(\mathbf{X}, \mathbf{Z}|\theta) d\mathbf{Z} = \\ &= \arg \max_{\theta} \log \int p(\mathbf{X}|\mathbf{Z}, \theta) p(\mathbf{Z}) d\mathbf{Z}.\end{aligned}$$

## Variational lower bound

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

## Kullback-Leibler divergence

- ▶  $KL(q||p) \geq 0$ ;
- ▶  $KL(q||p) = 0 \Leftrightarrow q \equiv p$ .



## Variational lower bound

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

### ELBO

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

Instead of maximizing incomplete likelihood, maximize ELBO

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

# 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.

# Amortized variational inference

## E-step

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

could be **intractable**.

## Idea

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

## Variational Bayes

- ▶ E-step

$$\phi_n = \phi_{n-1} + \eta \nabla_{\phi} \mathcal{L}(\phi, \theta_{n-1})|_{\phi=\phi_{n-1}}$$

- ▶ M-step

$$\theta_n = \theta_{n-1} + \eta \nabla_{\theta} \mathcal{L}(\phi_n, \theta)|_{\theta=\theta_{n-1}}$$

## ELBO gradient (M-step)

$$\mathcal{L}(\phi, \theta) = \mathbb{E}_q \log p(\mathbf{X}|\mathbf{Z}, \theta) - KL(q(\mathbf{Z}|\mathbf{X}, \phi)||p(\mathbf{Z})) \rightarrow \max_{\phi, \theta}.$$

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

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

## ELBO gradient (E-step)

$$\mathcal{L}(\phi, \theta) = \mathbb{E}_q \log p(\mathbf{X}|\mathbf{Z}, \theta) - KL(q(\mathbf{Z}|\mathbf{X}, \phi)||p(\mathbf{Z})) \rightarrow \max_{\phi, \theta}.$$

Optimization w.r.t.  $\phi$ : density function depends on the parameters.

Hint 1 (log-derivative trick)

$$\nabla_x p(y|x) = p(y|x) \nabla_x \log p(y|x).$$

Hint 2

$$\begin{aligned} \nabla_x f(x) &= \nabla_x \int p(y|x) h(y) dy \\ &= \int (\nabla_x p(y|x)) h(y) dy \\ &\sim h(y_0) \nabla_x \log p(y_0|x) \quad y_0 \sim p(y|x). \end{aligned}$$

## ELBO gradient (E-step)

$$\mathcal{L}(\phi, \theta) = \mathbb{E}_q \log p(\mathbf{X}|\mathbf{Z}, \theta) - KL(q(\mathbf{Z}|\mathbf{X}, \phi)||p(\mathbf{Z})) \rightarrow \max_{\phi, \theta}.$$

Optimization w.r.t.  $\phi$ : density function depends on the parameters.

$$\nabla_{\phi} \int q(\mathbf{Z}|\mathbf{X}, \phi) \log p(\mathbf{X}|\mathbf{Z}, \theta) d\mathbf{Z} \sim \log p(\mathbf{x}_i|\mathbf{z}_i^*, \theta) \nabla_{\phi} \log q(\mathbf{z}_i^*|\mathbf{x}_i, \phi),$$

$$\mathbf{z}_i^* \sim q(\mathbf{z}_i^*|\mathbf{x}_i, \phi).$$

### Problem

Unstable solution with huge variance.

### Solution

Reparametrization trick

# ELBO gradient (E-step)

Reparametrization trick

$$f(x) = \int p(y|x)h(y)dy$$

$$\begin{aligned}\nabla_x \int p(y|x)h(y)dy &= \nabla_x \int r(\epsilon)h(g(x, \epsilon))d\epsilon \\ &= \nabla_x h(g(x, \epsilon^*)), \quad \epsilon^* \sim r(\epsilon).\end{aligned}$$

Example

$$q(z|x) = \mathcal{N}(z|\mu, \sigma^2), \quad r(\epsilon) = \mathcal{N}(\epsilon|0, 1), \quad z = \sigma\epsilon + \mu.$$

# ELBO gradient (E-step)

## Derivative

$$\begin{aligned}\nabla_{\phi} \int q(\mathbf{Z}|\mathbf{X}, \phi) \log p(\mathbf{X}|\mathbf{Z}, \theta) d\mathbf{Z} &\sim \\ n \nabla_{\phi} \int r(\epsilon) \log p(\mathbf{x}_i | g(\mathbf{x}_i, \epsilon, \phi), \theta) d\epsilon &\sim \\ n \nabla_{\phi} \log p(\mathbf{x}_i | g(\mathbf{x}_i, \epsilon^*, \phi), \theta), &\quad \epsilon^* \sim r(\epsilon).\end{aligned}$$

## Variational assumption

$$q(\mathbf{z}|\mathbf{x}, \theta) = \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}), \boldsymbol{\Sigma}(\mathbf{x})).$$



# Variational autoencoder (VAE)

## Final algorithm

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

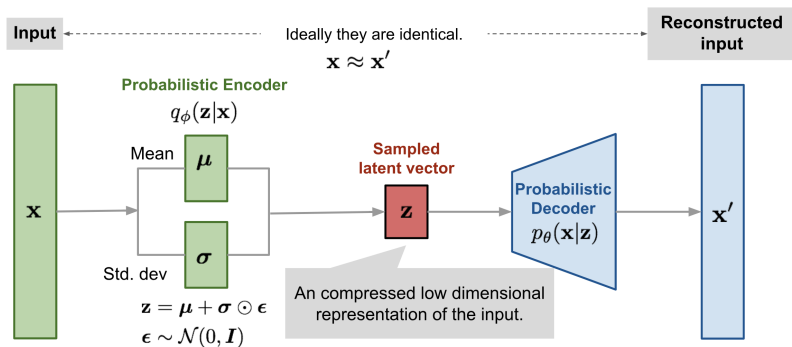
$$n \nabla_{\theta} \log p(\mathbf{x}_i | \mathbf{z}_i^*, \theta), \quad \mathbf{z}_i^* \sim q(\mathbf{z}_i | \mathbf{x}_i, \phi);$$

- ▶ compute stochastic gradient w.r.t.  $\phi$

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

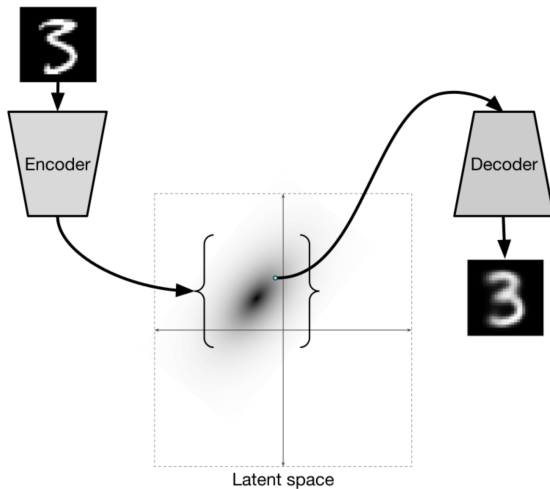
- ▶ update  $\theta, \phi$  according to the selected optimization method (SGD, Adam, RMSProp).

# Variational autoencoder (VAE)



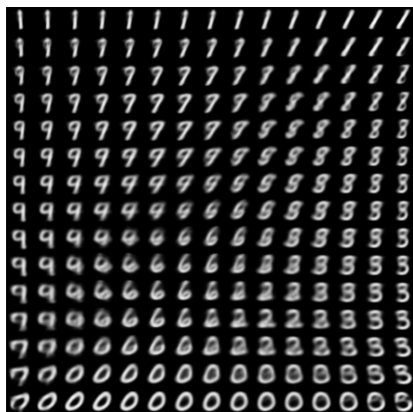
<https://lilianweng.github.io/lil-log/2018/08/12/from-autoencoder-to-beta-vaе.html>

# Variational Autoencoder



# Variational Autoencoder

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



<http://bit.ly/2w73aXB>

# References

- ▶ *Variational Bayesian inference with Stochastic Search*  
<https://arxiv.org/abs/1206.6430>
- ▶ *Stochastic Variational Inference*  
<https://arxiv.org/abs/1206.7051>
- ▶ *Doubly Stochastic Variational Bayes for non-Conjugate Inference*  
<http://proceedings.mlr.press/v32/titsias14.pdf>
- ▶ *Auto-Encoding Variational Bayes*  
<https://arxiv.org/abs/1312.6114>
- ▶ *Markov chain Monte Carlo and variational inference: Bridging the gap*  
<https://arxiv.org/pdf/1410.6460.pdf>
- ▶ *Tutorial on Variational Autoencoders*  
<http://arxiv.org/abs/1606.05908>