

# Deep Generative Models

Roman Isachenko

Moscow Institute of Physics and Technology

2019

# Likelihood-based models

## Exact likelihood evaluation

- ▶ Autoregressive models (PixelCNN, WaveNet);
- ▶ Flow models (NICE, RealNVP).

## Approximate likelihood evaluation

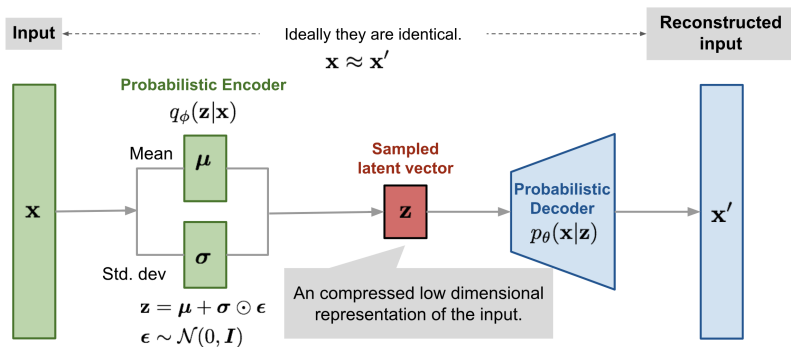
- ▶ Latent variable models (VAE).

What are the pros and cons of each of them?

How are they connected?

# VAE recap

$$p(\mathbf{x}|\theta) \geq \mathcal{L}(\phi, \theta) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \phi)} \log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x}, \phi)} \rightarrow \max_{\phi, \theta}.$$



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

## VAE limitations

- ▶ Poor variational posterior distribution (encoder)

$$q(\mathbf{z}|\mathbf{x}, \phi) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_\phi(\mathbf{x}), \boldsymbol{\sigma}_\phi^2(\mathbf{x})).$$

- ▶ Poor prior distribution

$$p(\mathbf{z}) = \mathcal{N}(0, \mathbf{I}).$$

- ▶ Poor probabilistic model (decoder)

$$p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_\theta(\mathbf{z}), \boldsymbol{\sigma}_\theta^2(\mathbf{z})).$$

- ▶ Loose lower bound

$$p(\mathbf{x}|\boldsymbol{\theta}) - \mathcal{L}(q, \boldsymbol{\theta}) = (?).$$

# Variational posterior

We wish  $KL(q(\mathbf{z}|\mathbf{x}, \phi) || p(\mathbf{z}|\mathbf{x}, \theta)) = 0$ .

(In this case the lower bound is tight  $p(\mathbf{x}|\theta) = \mathcal{L}(q, \theta)$ ).

Normal variational distribution  $q(\mathbf{z}|\mathbf{x}, \phi) = \mathcal{N}(\mathbf{z}|\mu_\phi(\mathbf{x}), \sigma_\phi^2(\mathbf{x}))$  is poor (e.g. has only one mode).

Flows models transform simple base distribution to complex one using invertible transformation with simple Jacobian.

How to use flows in VAE?

# Flows in VAE

Apply the sequence of transformations to the random variables

$$\mathbf{z}_0 \sim q_0(\mathbf{z}|\mathbf{x}, \phi) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_\phi(\mathbf{x}), \boldsymbol{\sigma}_\phi^2(\mathbf{x})).$$

Here,  $q_0(\mathbf{z}|\mathbf{x}, \phi)$  plays the role of a base distribution.

$$\mathbf{z}_0 \xrightarrow{g_1} \mathbf{z}_1 \xrightarrow{g_2} \dots \xrightarrow{g_K} \mathbf{z}_K.$$

Each  $g_k$  is a flow transformation (e.g. planar, radial, coupling layer).

$$\log q_K(\mathbf{z}_K) = \log q_0(\mathbf{z}_0) - \sum_{k=1}^K \log \left| \det \left( \frac{\partial g_k(\mathbf{z}_{k-1})}{\partial \mathbf{z}_{k-1}} \right) \right|.$$

---

<https://arxiv.org/pdf/1505.05770.pdf>

## Flows in VAE

$$\log q_K(\mathbf{z}_K) = \log q_0(\mathbf{z}_0) - \sum_{k=1}^K \log \left| \det \left( \frac{\partial \mathbf{g}_k(\mathbf{z}_{k-1})}{\partial \mathbf{z}_{k-1}} \right) \right|.$$

Now the variational posterior is  $q_K(\mathbf{z}_K|\mathbf{x}, \phi)$ .

$$\begin{aligned} \mathcal{L}(\phi, \theta) &= \mathbb{E}_{q_K(\mathbf{z}_K|\mathbf{x}, \phi)} \log \frac{p(\mathbf{x}, \mathbf{z}_K|\theta)}{q_K(\mathbf{z}_K|\mathbf{x}, \phi)} \\ &= \mathbb{E}_{q_K(\mathbf{z}_K|\mathbf{x}, \phi)} [\log p(\mathbf{x}, \mathbf{z}_K|\theta) - \log q_K(\mathbf{z}_K|\mathbf{x}, \phi)] \\ &= \mathbb{E}_{q_0(\mathbf{z}_0|\mathbf{x}, \phi)} [\log p(\mathbf{x}, \mathbf{z}_K|\theta) - \log q_K(\mathbf{z}_K|\mathbf{x}, \phi)] \\ &= \mathbb{E}_{q_0(\mathbf{z}_0|\mathbf{x}, \phi)} \left[ \log p(\mathbf{x}, \mathbf{z}_K|\theta) - \log q_0(\mathbf{z}_0|\mathbf{x}, \phi) - \right. \\ &\quad \left. - \sum_{k=1}^K \log \left| \det \left( \frac{\partial \mathbf{g}_k(\mathbf{z}_{k-1})}{\partial \mathbf{z}_{k-1}} \right) \right| \right]. \end{aligned}$$

# MAF, 2017

Consider autoregressive model

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

with conditionals

$$p(x_i|\mathbf{x}_{1:i-1}, \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}_{i,\boldsymbol{\theta}}(\mathbf{x}_{1:i-1}), \boldsymbol{\sigma}_{i,\boldsymbol{\theta}}^2(\mathbf{x}_{1:i-1})).$$

## Sampling

$$x_i = \boldsymbol{\sigma}_{i,\boldsymbol{\theta}}(\mathbf{x}_{1:i-1}) \cdot z_i + \boldsymbol{\mu}_{i,\boldsymbol{\theta}}(\mathbf{x}_{1:i-1}), \quad z_i \sim \mathcal{N}(0, 1).$$

## Inverse transform

$$z_i = (x_i - \boldsymbol{\mu}_{i,\boldsymbol{\theta}}(\mathbf{x}_{1:i-1})) \cdot \frac{1}{\boldsymbol{\sigma}_{i,\boldsymbol{\theta}}(\mathbf{x}_{1:i-1})}.$$



## Sampling

$$\mathbf{x} = g(\mathbf{z}, \theta) = \sigma_{\theta}(\mathbf{x}) \odot \mathbf{z} + \mu_{\theta}(\mathbf{x}), \quad \mathbf{z} \sim \mathcal{N}(0, \mathbf{I}).$$

## Inverse transform

$$\mathbf{z} = f(\mathbf{x}, \theta) = (\mathbf{x} - \mu_{\theta}(\mathbf{x})) \odot \frac{1}{\sigma_{\theta}(\mathbf{x})}.$$

What is the Jacobian of such flow?

- ▶ Sampling is slow (sequential).
- ▶ Likelihood evaluation is fast (e.g. MADE).

Suitable for density evaluation task.

---

<https://arxiv.org/pdf/1705.07057.pdf>

## Inverse transform in MAF

$$\begin{aligned}
 z_i &= (x_i - \boldsymbol{\mu}_{i,\boldsymbol{\theta}}(\mathbf{x}_{1:i-1})) \cdot \frac{1}{\boldsymbol{\sigma}_{i,\boldsymbol{\theta}}(\mathbf{x}_{1:i-1})} \\
 &= \frac{x_i}{\boldsymbol{\sigma}_{i,\boldsymbol{\theta}}(\mathbf{x}_{1:i-1})} - \frac{\boldsymbol{\mu}_{i,\boldsymbol{\theta}}(\mathbf{x}_{1:i-1})}{\boldsymbol{\sigma}_{i,\boldsymbol{\theta}}(\mathbf{x}_{1:i-1})} \\
 &= \hat{\boldsymbol{\sigma}}_{i,\boldsymbol{\theta}}(\mathbf{x}_{1:i-1}) \cdot x_i + \hat{\boldsymbol{\mu}}_{i,\boldsymbol{\theta}}(\mathbf{x}_{1:i-1}).
 \end{aligned}$$

How to make this transform to be efficient for sampling?

---

<https://arxiv.org/pdf/1606.04934.pdf>

## Sampling and inverse transform in MAF

$$x_i = \sigma_{i,\theta}(\mathbf{x}_{1:i-1}) \cdot z_i + \mu_{i,\theta}(\mathbf{x}_{1:i-1}).$$

$$z_i = (x_i - \mu_{i,\theta}(\mathbf{x}_{1:i-1})) \cdot \frac{1}{\sigma_{i,\theta}(\mathbf{x}_{1:i-1})}.$$

## Sampling and inverse transform in IAF

$$x_i = \hat{\sigma}_{i,\theta}(\mathbf{z}_{1:i-1}) \cdot z_i + \hat{\mu}_{i,\theta}(\mathbf{z}_{1:i-1}).$$

$$z_i = (x_i - \hat{\mu}_{i,\theta}(\mathbf{z}_{1:i-1})) \cdot \frac{1}{\hat{\sigma}_{i,\theta}(\mathbf{z}_{1:i-1})}.$$

MAF and IAF are inverse to each other up to reparametrization

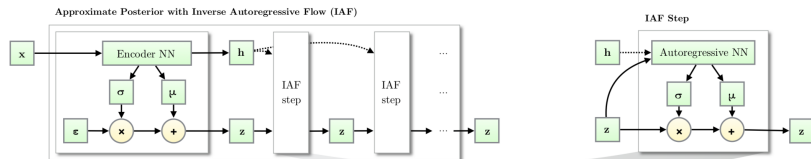
$$\hat{\sigma}_i = \frac{1}{\sigma_i}; \quad \hat{\mu}_i = \frac{\mu_i}{\sigma_i}.$$

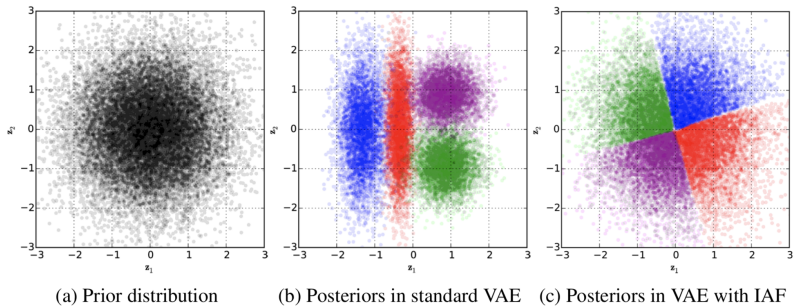
## Sampling

$$x_i = \hat{\sigma}_{i,\theta}(\mathbf{z}_{1:i-1}) \cdot z_i + \hat{\mu}_{i,\theta}(\mathbf{z}_{1:i-1}).$$

## Approximate posterior

$$z_i = (x_i - \hat{\mu}_{i,\theta}(\mathbf{z}_{1:i-1})) \cdot \frac{1}{\hat{\sigma}_{i,\theta}(\mathbf{z}_{1:i-1})}.$$





<https://arxiv.org/pdf/1606.04934.pdf>

# MAF vs IAF

## Theorem

Training a MAF with maximum likelihood corresponds to fitting an implicit IAF to the base density with stochastic variational inference:

$$\max_{\theta} p(\mathbf{X}|\theta) \Leftrightarrow \min_{\theta} KL(p(\mathbf{z}|\theta)||\pi(\mathbf{z}))$$

(Here,  $\pi(\mathbf{z})$  is a base distribution,  $\pi(\mathbf{x})$  is a data distribution).

## Proof

$$\begin{aligned} KL(p(\mathbf{z}|\theta)||\pi(\mathbf{z})) &= \mathbb{E}_{p(\mathbf{z}|\theta)} [\log p(\mathbf{z}|\theta) - \log \pi(\mathbf{z})] = \\ &= \mathbb{E}_{p(\mathbf{z}|\theta)} \left[ \log \pi(g(\mathbf{z})) + \log \left| \det \left( \frac{\partial g(\mathbf{z})}{\partial \mathbf{z}} \right) \right| - \log \pi(\mathbf{z}) \right] = \\ &= \mathbb{E}_{\pi(\mathbf{x})} \left[ \log \pi(\mathbf{x}) - \log \left| \det \left( \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \right) \right| - \log \pi(f(\mathbf{x})) \right]. \end{aligned}$$

# MAF vs IAF

## Proof (continued)

$$\begin{aligned} KL(p(\mathbf{z}|\boldsymbol{\theta})||\pi(\mathbf{z})) &= \mathbb{E}_{\pi(\mathbf{x})} \left[ \log \pi(\mathbf{x}) - \log \left| \det \left( \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \right) \right| - \log \pi(f(\mathbf{x})) \right] = \\ &= \mathbb{E}_{\pi(\mathbf{x})} [\log \pi(\mathbf{x}) - \log p(\mathbf{x}|\boldsymbol{\theta})] = KL(\pi(\mathbf{x})||p(\mathbf{x}|\boldsymbol{\theta})). \end{aligned}$$

$$\begin{aligned} \arg \min_{\boldsymbol{\theta}} KL(\pi(\mathbf{x})||p(\mathbf{x}|\boldsymbol{\theta})) &= \arg \min_{\boldsymbol{\theta}} \mathbb{E}_{\pi(\mathbf{x})} [\log \pi(\mathbf{x}) - \log p(\mathbf{x}|\boldsymbol{\theta})] \\ &= \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{\pi(\mathbf{x})} \log p(\mathbf{x}|\boldsymbol{\theta}) \end{aligned}$$

Unbiased estimator is MLE:

$$\mathbb{E}_{\pi(\mathbf{x})} \log p(\mathbf{x}|\boldsymbol{\theta}) = \sum_{i=1}^n \log p(\mathbf{x}_i|\boldsymbol{\theta}).$$

# MAF vs IAF vs RealNVP

## RealNVP

$$\begin{aligned}\mathbf{x}_{1:d} &= \mathbf{z}_{1:d}; \\ \mathbf{x}_{d:m} &= \mathbf{z}_{d:m} \odot \exp(c_1(\mathbf{z}_{1:d}, \boldsymbol{\theta})) + c_z(\mathbf{x}_{1:d}, \boldsymbol{\theta})\end{aligned}$$

## MAF

$$\mathbf{x} = \boldsymbol{\sigma}_\theta(\mathbf{x}) \odot \mathbf{z} + \boldsymbol{\mu}_\theta(\mathbf{x}).$$

## IAF

$$\mathbf{x} = \hat{\boldsymbol{\sigma}}_\theta(\mathbf{z}) \odot \mathbf{z} + \hat{\boldsymbol{\mu}}_\theta(\mathbf{z}).$$

How they are connected? Which flow is the most flexible?

---

<https://arxiv.org/pdf/1705.07057.pdf>



# MAF/IAF pros and cons

## MAF

- ▶ Sampling is slow.
- ▶ Likelihood evaluation is fast.

## IAF

- ▶ Sampling is fast.
- ▶ Likelihood evaluation is slow.

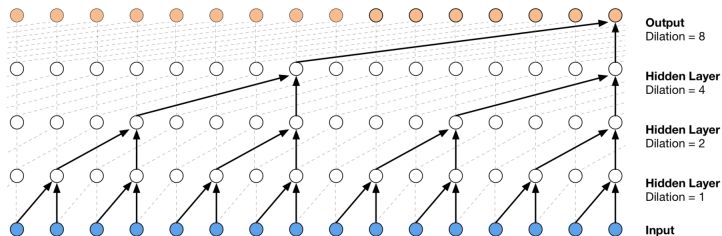
How to take the best of both worlds?

# WaveNet (2016)

Autoregressive model for raw audio waveforms generation

$$p(\mathbf{x}|\theta) = \prod_{t=1}^T p(x_t|\mathbf{x}_{1:t-1}, \theta).$$

The model uses causal dilated convolutions.



<https://arxiv.org/pdf/1609.03499.pdf>

# Parallel WaveNet, 2017

## Previous WaveNet model

- ▶ raw audio is high-dimensional (e.g. 16000 samples per second for 16kHz audio);
- ▶ WaveNet encodes 8-bit signal with 256-way categorical distribution.

## Goal

- ▶ improved fidelity (24kHz instead of 16kHz) → increase dilated convolution filter size from 2 to 3;
- ▶ 16-bit signals → mixture of logistics instead of categorical distribution.

---

<https://arxiv.org/pdf/1711.10433.pdf>

# Parallel WaveNet, 2017

## Probability density distillation

1. Train usual WaveNet (MAF) via MLE (teacher network).
2. Train IAF WaveNet model (student network), which attempts to match the probability of its own samples under the distribution learned by the teacher.

## Student objective

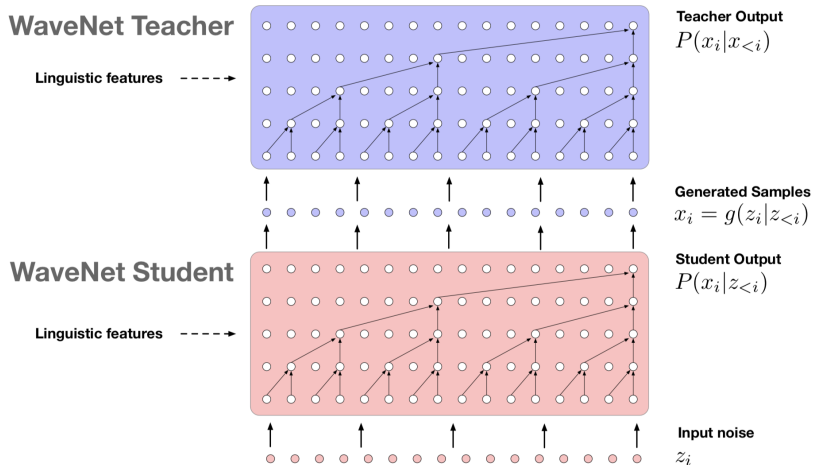
$$KL(p_s || p_t) = H(p_s, p_t) - H(p_s).$$

More than 1000x speed-up relative to original WaveNet!

---

<https://arxiv.org/pdf/1711.10433.pdf>

# Parallel WaveNet, 2017



<https://arxiv.org/pdf/1711.10433.pdf>

# References

- ▶ **NICE: Non-Independent Components Estimation**  
<https://arxiv.org/abs/1410.8516>  
**Summary:** Uses flows to model complex high-dimensional densities. Introduce the ways to compute determinant of Jacobian in a simple way. Triangular Jacobian, coupling layers, factorized distribution.
- ▶ **Variational Inference with Normalizing Flows**  
<https://arxiv.org/abs/1505.05770>  
**Summary:** Propose to use normalizing flows in variational inference. Discuss finite and infinitesimal flows. Useful invertible flows: planar, radial. Comparison with NICE.
- ▶ **RealNVP: Density estimation using Real NVP**  
<https://arxiv.org/pdf/1605.08803.pdf>  
**Summary:** Authors of NICE. The same idea and architecture, more practical. Lots of experiments and images. Coupling layers with checkerboard and channel-wise permutations.
- ▶ **IAF: Improving Variational Inference with Inverse Autoregressive Flow**  
<https://arxiv.org/abs/1606.04934>  
**Summary:** Introduce inverse autoregressive flow (IAF). Models each autoregressive conditional as gaussian with autoregressive means and covariances. Inverse transformation allows to parallelize sampling.
- ▶ **MAF: Masked Autoregressive Flow for Density Estimation**  
<https://arxiv.org/pdf/1705.07057.pdf>  
**Summary:** Similar to IAF. Give comprehensive overview with link to IAF and RealNVP. MAF is suitable for density estimation, IAF as a recognition network.
- ▶ **Parallel WaveNet: Fast High-Fidelity Speech Synthesis**  
<https://arxiv.org/pdf/1711.10433.pdf>  
**Summary:** WaveNet is MAF (sequential generation). To exploit IAF fast sampling, knowledge distillation used. Teacher network is large WaveNet, student - is a IAF small WaveNet (generate samples from noise is parallel). The loss is KL divergence between student and teacher distributions. The additional perceptual, contrastive and power losses used to create more natural sounds.