# Deep Generative Models

## Lecture 6

Roman Isachenko

Moscow Institute of Physics and Technology

Autumn, 2021

# Recap of previous lecture



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

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

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

## Flow likelihood

$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \log p(f(\mathbf{x}, \boldsymbol{\theta})) + \log \left| \det \left( \frac{\partial f(\mathbf{x}, \boldsymbol{\theta})}{\partial \mathbf{x}} \right) \right|$$

## What we want

- Efficient computation of Jacobian $\frac{\partial f(\mathbf{x}, \boldsymbol{\theta})}{\partial \mathbf{x}}$;
- Efficient sampling from the base distribution $p(\mathbf{z})$;
- Efficient inversion of $f(\mathbf{x}, \boldsymbol{\theta})$.

Dinh L., Sohl-Dickstein J., Bengio S. Density estimation using Real NVP, 2016

# Recap of previous lecture

## Planar flow

$$g(\mathbf{z}, \boldsymbol{\theta}) = \mathbf{z} + \mathbf{u}\,h(\mathbf{w}^T\mathbf{z} + b).$$

## Sylvester flow

$$g(\mathbf{z}, \boldsymbol{\theta}) = \mathbf{z} + \mathbf{A}\,h(\mathbf{B}\mathbf{z} + \mathbf{b}).$$

## NICE/RealNVP: Affine coupling law

$$\begin{cases} \mathbf{z}_{1:d} = \mathbf{x}_{1:d}; \\ \mathbf{z}_{d:m} = \mathbf{x}_{d:m} \odot \exp\left(c_1(\mathbf{x}_{1:d}, \boldsymbol{\theta})\right) + c_2(\mathbf{x}_{1:d}, \boldsymbol{\theta}). \end{cases}$$

## Glow: invertible 1x1 conv

$$\mathbf{W} = \mathbf{P}\mathbf{L}(\mathbf{U} + \mathrm{diag}(\mathbf{s})).$$

Rezende D. J., Mohamed S. Variational Inference with Normalizing Flows, 2015
Berg R. et al. Sylvester normalizing flows for variational inference, 2018
Dinh L., Krueger D., Bengio Y. NICE: Non-linear Independent Components Estimation, 2014
Dinh L., Sohl-Dickstein J., Bengio S. Density estimation using Real NVP, 2016
Kingma D. P., Dhariwal P. Glow: Generative Flow with Invertible 1x1 Convolutions, 2018

# Likelihood-based models

### Exact likelihood evaluation

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

### Approximate likelihood evaluation

- ▶ Latent variable models (VAE).

What are the pros and cons of each of them?

# VAE recap

### ELBO

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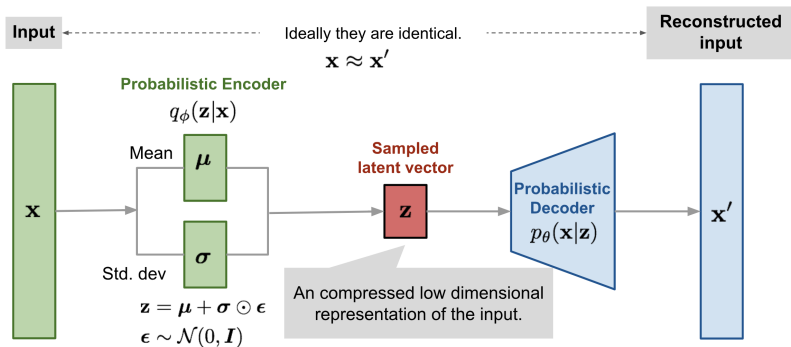


image credit:
https://lilianweng.github.io/lil-log/2018/08/12/from-autoencoder-to-beta-vae.html

# VAE limitations

▶ Poor variational posterior distribution (encoder)

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

▶ Poor prior distribution

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

▶ Poor probabilistic model (decoder)

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

▶ Loose lower bound

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

# Variational posterior

### ELBO

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

▶ In E-step of EM-algorithm we wish
$KL(q(\mathbf{z}|\mathbf{x}, \boldsymbol{\phi})||p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})) = 0$.
(In this case the lower bound is tight $\log p(\mathbf{x}|\boldsymbol{\theta}) = \mathcal{L}(q, \boldsymbol{\theta})$).

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

▶ Flows models convert a simple base distribution to a compex one using invertible transformation with simple Jacobian. How to use flows in VAE?

## Flows in VAE

Apply a sequence of transformations to the random variable

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

Here, $q(\mathbf{z}|\mathbf{x}, \phi)$ (which is a VAE encoder) plays a role of a base distribution.

$$\mathbf{z}_0 \xrightarrow{g_1} \mathbf{z}_1 \xrightarrow{g_2} \dots \xrightarrow{g_K} \mathbf{z}_K, \quad \mathbf{z}_K = g(\mathbf{z}_0), \quad g = g_K \circ \cdots \circ g_1.$$

Each $g_k$ is a flow transformation (e.g. planar, coupling layer) parameterized by $\phi_k$.

$$\begin{aligned}
\log q_K(\mathbf{z}_K|\mathbf{x}, \phi, \{\phi_k\}_{k=1}^K) = \log q(\mathbf{z}_0|\mathbf{x}, \phi) \\
- \sum_{k=1}^K \log \left| \det \left( \frac{\partial g_k(\mathbf{z}_{k-1}, \phi_k)}{\partial \mathbf{z}_{k-1}} \right) \right|.
\end{aligned}$$

---

*Rezende D. J., Mohamed S. Variational Inference with Normalizing Flows, 2015*

# Flows in VAE

## ELBO

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

## Flow model in latent space

$$\log q_K(\mathbf{z}_K|\mathbf{x}, \boldsymbol{\phi}_*) = \log q(\mathbf{z}_0|\mathbf{x}, \boldsymbol{\phi}) - \sum_{k=1}^{K} \log \left| \det \left( \frac{\partial g_k(\mathbf{z}_{k-1}, \boldsymbol{\phi}_k)}{\partial \mathbf{z}_{k-1}} \right) \right|.$$

Let use $q_K(\mathbf{z}_K|\mathbf{x}, \boldsymbol{\phi}_*)$, $\boldsymbol{\phi}_* = \{\boldsymbol{\phi}, \boldsymbol{\phi}_1, \ldots, \boldsymbol{\phi}_K\}$ as a variational distribution. Here $\boldsymbol{\phi}$ – encoder parameters, $\{\boldsymbol{\phi}_k\}_{k=1}^{K}$ – flow parameters.

- Encoder outputs base distribution $q(\mathbf{z}_0|\mathbf{x}, \boldsymbol{\phi})$.
- Flow model $\mathbf{z}_K = g(\mathbf{z}_0, \{\boldsymbol{\phi}_k\}_{k=1}^{K})$ transforms the base distribution $q(\mathbf{z}_0|\mathbf{x}, \boldsymbol{\phi})$ to the distribution $q_K(\mathbf{z}_K|\mathbf{x}, \boldsymbol{\phi}_*)$.
- Distribution $q_K(\mathbf{z}_K|\mathbf{x}, \boldsymbol{\phi}_*)$ is used as a variational distribution for ELBO maximization.

---

Rezende D. J., Mohamed S. Variational Inference with Normalizing Flows, 2015

# Flows in VAE

### Flow model in latent space

$$\log q_K(\mathbf{z}_K|\mathbf{x}, \phi_*) = \log q(\mathbf{z}_0|\mathbf{x}, \phi) - \sum_{k=1}^{K} \log \left| \det \left( \frac{\partial g_k(\mathbf{z}_{k-1}, \phi_k)}{\partial \mathbf{z}_{k-1}} \right) \right|.$$

### ELBO objective

$$\begin{aligned}
\mathcal{L}(\phi, \boldsymbol{\theta}) &= \mathbb{E}_{q_K(\mathbf{z}_K|\mathbf{x}, \phi_*)} \log \frac{p(\mathbf{x}, \mathbf{z}_K|\boldsymbol{\theta})}{q_K(\mathbf{z}_K|\mathbf{x}, \phi_*)} \\
&= \mathbb{E}_{q_K(\mathbf{z}_K|\mathbf{x}, \phi_*)} \big[ \log p(\mathbf{x}, \mathbf{z}_K|\boldsymbol{\theta}) - \log q_K(\mathbf{z}_K|\mathbf{x}, \phi_*) \big] \\
&= \mathbb{E}_{q_K(\mathbf{z}_K|\mathbf{x}, \phi_*)} \log p(\mathbf{x}|\mathbf{z}_K, \boldsymbol{\theta}) - KL(q_K(\mathbf{z}_K|\mathbf{x}, \phi_*)||p(\mathbf{z}_K)).
\end{aligned}$$

The second term in ELBO is reverse KL divergence. Planar flows was originally proposed for variational inference in VAE.

---

*Rezende D. J., Mohamed S. Variational Inference with Normalizing Flows, 2015*

# Flows in VAE

### Variational distribution

$$\log q_K(\mathbf{z}_K|\mathbf{x}, \phi_*) = \log q(\mathbf{z}_0|\mathbf{x}, \phi) - \sum_{k=1}^{K} \log \left| \det \left( \frac{\partial g_k(\mathbf{z}_{k-1}, \phi_k)}{\partial \mathbf{z}_{k-1}} \right) \right|.$$

### ELBO objective

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

Rezende D. J., Mohamed S. Variational Inference with Normalizing Flows, 2015

# Flows in VAE

### Variational distribution

$$\log q_K(\mathbf{z}_K|\mathbf{x}, \phi_*) = \log q(\mathbf{z}_0|\mathbf{x}, \phi) - \sum_{k=1}^{K} \log \left| \det \left( \frac{\partial g_k(\mathbf{z}_{k-1}, \phi_k)}{\partial \mathbf{z}_{k-1}} \right) \right|.$$

### ELBO objective

$$\mathcal{L}(\phi, \theta) = \mathbb{E}_{q(\mathbf{z}_0|\mathbf{x}, \phi)} \Bigg[ \log p(\mathbf{x}, \mathbf{z}_K|\theta) - \log q(\mathbf{z}_0|\mathbf{x}, \phi) +$$

$$+ \sum_{k=1}^{K} \log \left| \det \left( \frac{\partial g_k(\mathbf{z}_{k-1}, \phi_k)}{\partial \mathbf{z}_{k-1}} \right) \right| \Bigg].$$

- ▶ Obtain samples $\mathbf{z}_0$ from the encoder.
- ▶ Apply flow model $\mathbf{z}_K = g(\mathbf{z}_0, \{\phi_k\}_{k=1}^{K})$.
- ▶ Compute likelihood for $\mathbf{z}_K$ using the decoder, base distribution for $\mathbf{z}_0$ and the Jacobian.
- ▶ We do not need inverse flow function, if we use flows in variational inference.

Rezende D. J., Mohamed S. Variational Inference with Normalizing Flows, 2015

# Recap of previous lecture

### ELBO

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

- Normal variational distribution
  $q(\mathbf{z}|\mathbf{x}, \boldsymbol{\phi}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_{\boldsymbol{\phi}}(\mathbf{x}), \boldsymbol{\sigma}^2_{\boldsymbol{\phi}}(\mathbf{x}))$ is poor (e.g. has only one mode).
- Flows models convert a simple base distribution to a compex one using an invertible transformation with simple Jacobian.

### Flow model in latent space

$$\log q_K(\mathbf{z}_K|\mathbf{x}, \boldsymbol{\phi}_*) = \log q(\mathbf{z}_0|\mathbf{x}, \boldsymbol{\phi}) - \sum_{k=1}^{K} \log \left| \det \left( \frac{\partial g_k(\mathbf{z}_{k-1}, \boldsymbol{\phi}_k)}{\partial \mathbf{z}_{k-1}} \right) \right|.$$

Let's use $q_K(\mathbf{z}_K|\mathbf{x}, \boldsymbol{\phi}_*)$, $\boldsymbol{\phi}_* = \{\boldsymbol{\phi}, \boldsymbol{\phi}_1, \ldots, \boldsymbol{\phi}_K\}$ as a variational distribution. Here, $\boldsymbol{\phi}$ – encoder parameters, $\{\boldsymbol{\phi}_k\}_{k=1}^{K}$ – flow parameters.

---

Rezende D. J., Mohamed S. Variational Inference with Normalizing Flows, 2015

# Recap of previous lecture

### Variational distribution

$$\log q_K(\mathbf{z}_K|\mathbf{x}, \phi_*) = \log q(\mathbf{z}_0|\mathbf{x}, \phi) - \sum_{k=1}^{K} \log \left| \det \left( \frac{\partial g_k(\mathbf{z}_{k-1}, \phi_k)}{\partial \mathbf{z}_{k-1}} \right) \right|.$$

### ELBO objective

$$\mathcal{L}(\phi, \boldsymbol{\theta}) = \mathbb{E}_{q(\mathbf{z}_0|\mathbf{x}, \phi)} \Bigg[ \log p(\mathbf{x}, \mathbf{z}_K|\boldsymbol{\theta}) - \log q(\mathbf{z}_0|\mathbf{x}, \phi) +$$
$$+ \sum_{k=1}^{K} \log \left| \det \left( \frac{\partial g_k(\mathbf{z}_{k-1}, \phi_k)}{\partial \mathbf{z}_{k-1}} \right) \right| \Bigg].$$

- ▶ Obtain samples $\mathbf{z}_0$ from the encoder.
- ▶ Apply flow model $\mathbf{z}_K = g(\mathbf{z}_0, \{\phi_k\}_{k=1}^{K})$.
- ▶ Compute likelihood for $\mathbf{z}_K$ using the decoder, base distribution for $\mathbf{z}_0$ and the Jacobian.
- ▶ We do not need an inverse flow function if we use flows in variational inference.

Rezende D. J., Mohamed S. Variational Inference with Normalizing Flows, 2015

# Inverse autoregressive flow (IAF)

$$\mathbf{x} = g(\mathbf{z}, \boldsymbol{\theta}) \quad \Rightarrow \quad x_i = \tilde{\sigma}_i(\mathbf{z}_{1:i-1}) \cdot z_i + \tilde{\mu}_i(\mathbf{z}_{1:i-1}).$$

$$\mathbf{z} = f(\mathbf{x}, \boldsymbol{\theta}) \quad \Rightarrow \quad z_i = (x_i - \tilde{\mu}_i(\mathbf{z}_{1:i-1})) \cdot \frac{1}{\tilde{\sigma}_i(\mathbf{z}_{1:i-1})}.$$
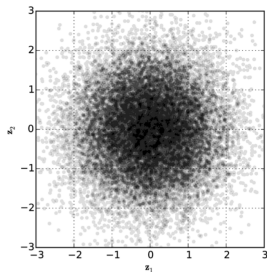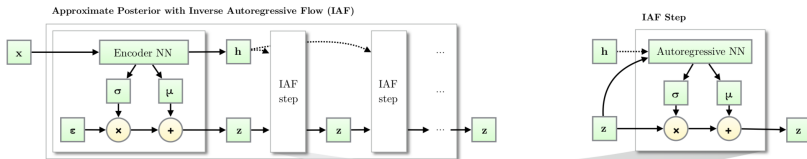
Reverse KL for flow model

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

- We don't need to think about computing the function $f(\mathbf{x}, \boldsymbol{\theta})$.
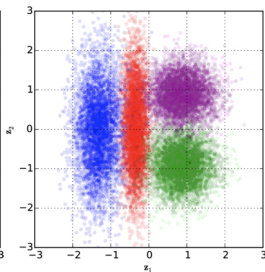- Inverse autoregressive flow is a natural choice for using flows in VAE:

$$\mathbf{z}_0 = \boldsymbol{\sigma}(\mathbf{x}) \odot \boldsymbol{\epsilon} + \boldsymbol{\mu}(\mathbf{x}), \quad \boldsymbol{\epsilon} \sim \mathcal{N}(0, 1); \quad \sim q(\mathbf{z}_0|\mathbf{x}, \phi).$$

$$\mathbf{z}_k = \tilde{\boldsymbol{\sigma}}_k(\mathbf{z}_{k-1}) \odot \mathbf{z}_{k-1} + \tilde{\boldsymbol{\mu}}_k(\mathbf{z}_{k-1}), \quad k \geq 1; \quad \sim q_k(\mathbf{z}_k|\mathbf{x}, \phi, \{\phi_j\}_{j=1}^k).$$
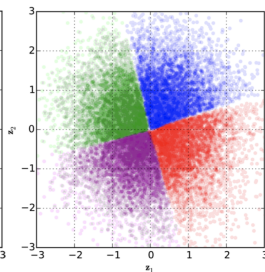
Kingma D. P. et al. *Improving Variational Inference with Inverse Autoregressive Flow*, 2016

# Inverse autoregressive flow (IAF)



(a) Prior distribution    (b) Posteriors in standard VAE    (c) Posteriors in VAE with IAF

*Kingma D. P. et al. Improving Variational Inference with Inverse Autoregressive Flow, 2016*

# MAF/IAF pros and cons

## MAF

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

## IAF

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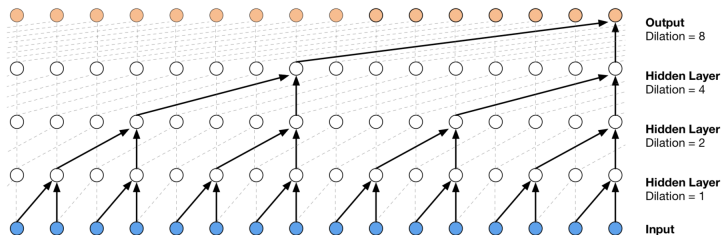- ► Likelihood evaluation is slow.

How to take the best of both worlds?

# WaveNet (2016)

Autoregressive model for raw audio waveforms generation

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

The model uses causal dilated convolutions.



Oord A. et al. Wavenet: A generative model for raw audio, 2016

# 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.

---

*Oord A. et al. Parallel WaveNet: Fast High-Fidelity Speech Synthesis, 2017*

# Parallel WaveNet, 2017
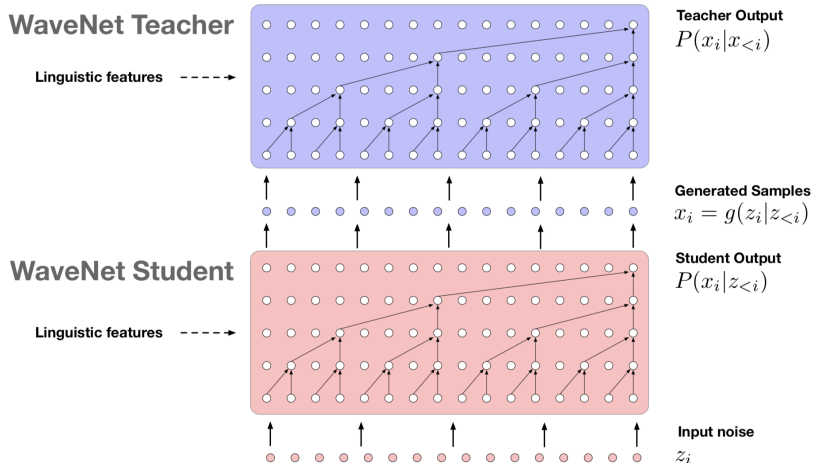
### Probability density distillation

1. Train usual WaveNet (MAF) via MLE (teacher network).
2. Train IAF WaveNet (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!

*Oord A. et al. Parallel WaveNet: Fast High-Fidelity Speech Synthesis, 2017*

# Parallel WaveNet, 2017



Oord A. et al. Parallel WaveNet: Fast High-Fidelity Speech Synthesis, 2017

# Flow KL duality

### Theorem

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

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

- $p(\mathbf{z})$ is a base distribution; $\pi(\mathbf{x})$ is a data distribution;
- $\mathbf{z} \sim p(\mathbf{z})$, $\mathbf{x} = g(\mathbf{z}, \boldsymbol{\theta})$, $\mathbf{x} \sim p(\mathbf{x}|\boldsymbol{\theta})$;
- $\mathbf{x} \sim \pi(\mathbf{x})$, $\mathbf{z} = f(\mathbf{x}, \boldsymbol{\theta})$, $\mathbf{z} \sim p(\mathbf{z}|\boldsymbol{\theta})$;

$$\log p(\mathbf{z}|\boldsymbol{\theta}) = \log \pi(g(\mathbf{z}, \boldsymbol{\theta})) + \log \left| \det \left( \frac{\partial g(\mathbf{z}, \boldsymbol{\theta})}{\partial \mathbf{z}} \right) \right|;$$

$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \log p(f(\mathbf{x}, \boldsymbol{\theta})) + \log \left| \det \left( \frac{\partial f(\mathbf{x}, \boldsymbol{\theta})}{\partial \mathbf{x}} \right) \right|.$$

Papamakarios G. et al. Normalizing flows for probabilistic modeling and inference, 2019

# MAF vs IAF

### Theorem

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

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

### Proof

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

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

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

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

# MAF vs IAF

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

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

### Proof (continued)

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

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

# Dequantization

- Images are discrete data, pixels lie in the [0, 255] integer domain (the model is $P(\mathbf{x}|\boldsymbol{\theta}) = \text{Categorical}(\boldsymbol{\pi}(\boldsymbol{\theta}))$).
- Flow is a continuous model (it works with continuous data $\mathbf{x}$).

By fitting a continuous density model to discrete data, one can produce a degenerate solution with all probability mass on discrete values.
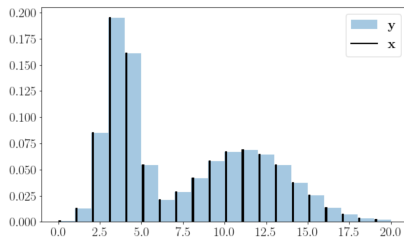
How to convert a discrete data distribution to a continuous one?

## Uniform dequantization

$$\mathbf{x} \sim \text{Categorical}(\boldsymbol{\pi})$$
$$\mathbf{u} \sim U[0, 1]$$

$$\mathbf{y} = \mathbf{x} + \mathbf{u} \sim \text{Continuous}$$



*Theis L., Oord A., Bethge M. A note on the evaluation of generative models. 2015*

# Uniform dequantization

### Statement

Fitting continuous model $p(\mathbf{y}|\boldsymbol{\theta})$ on uniformly dequantized data $\mathbf{y} = \mathbf{x} + \mathbf{u}$, $\mathbf{u} \sim U[0,1]$ is equivalent to maximization of a lower bound on log-likelihood for a discrete model:

$$P(\mathbf{x}|\boldsymbol{\theta}) = \int_{U[0,1]} p(\mathbf{x} + \mathbf{u}|\boldsymbol{\theta}) d\mathbf{u}$$

Thus, the maximisation of continuous model log-likelihood on $\mathbf{y}$ can't lead to the a collapse onto the discrete data (the objective is bounded above by the discrete model log-likelihood).

### Proof

$$\log P(\mathbf{x}|\boldsymbol{\theta}) = \log \int_{U[0,1]} p(\mathbf{x} + \mathbf{u}|\boldsymbol{\theta}) d\mathbf{u} \geq$$
$$\geq \int_{U[0,1]} \log p(\mathbf{x} + \mathbf{u}|\boldsymbol{\theta}) d\mathbf{u} = \log p(\mathbf{y}|\boldsymbol{\theta}).$$

# Summary

- Gaussian autoregressive model is a special type of flow.

- MAF is an example of such a model which is suitable for density estimation tasks. IAF uses an inverse autoregressive transformation for variational inference task.

- RealNVP is a special case of IAF and MAF.

- There is a duality between forward and reverse KL for flow models.

- To apply a continuous model to a discrete distribution it is standard practice to dequantize data at first.