

# Deep Generative Models

## Lecture 2

Roman Isachenko



2025, Spring

## Recap of previous lecture

We are given i.i.d. samples  $\{\mathbf{x}_i\}_{i=1}^n \in \mathbb{R}^m$  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  $\mathbf{x} \sim \pi(\mathbf{x})$ ).

Instead of searching true  $\pi(\mathbf{x})$  over all probability distributions, learn function approximation  $p(\mathbf{x}|\theta)$   $\approx \pi(\mathbf{x})$ .

### Divergence

- ▶  $D(\pi||p) \geq 0$  for all  $\pi, p \in \mathcal{P}$ ;
- ▶  $D(\pi||p) = 0$  if and only if  $\pi \equiv p$ .

### Divergence minimization task

$$\min_{\theta} D(\pi||p).$$

## Recap of previous lecture

Forward KL

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

Reverse KL

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

Maximum likelihood estimation (MLE)

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

Maximum likelihood estimation is equivalent to minimization of the Monte-Carlo estimate of forward KL.

## Recap of previous lecture

Likelihood as product of conditionals

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

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

MLE problem for autoregressive model

$$\theta^* = \arg \max_{\theta} \sum_{i=1}^n \sum_{j=1}^m \log p(\underline{x_{ij}}|\underline{\mathbf{x}_{i,1:j-1}}, \theta).$$

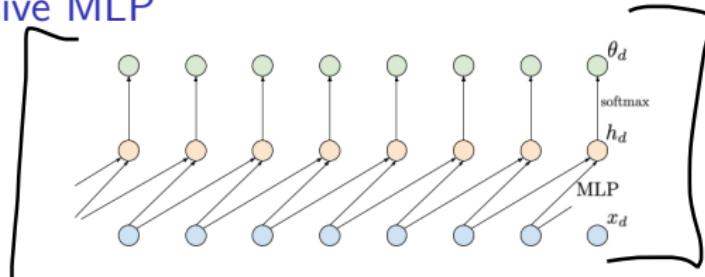
Sampling

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

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

# Recap of previous lecture

## Autoregressive MLP



## Autoregressive transformer

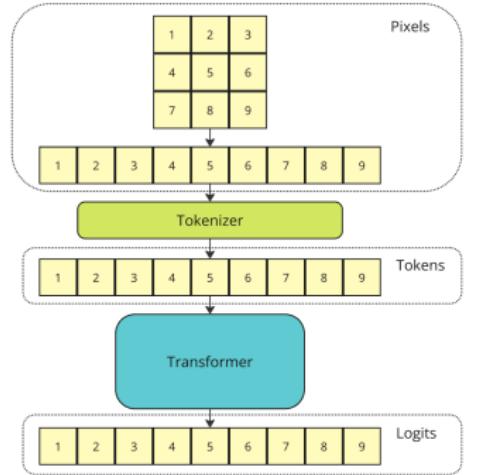
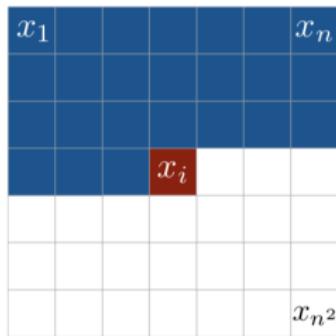


image credit: [https://jmtomczak.github.io/blog/2/2\\_ARM.html](https://jmtomczak.github.io/blog/2/2_ARM.html)

Chen M. et al. Generative Pretraining from Pixels, 2020

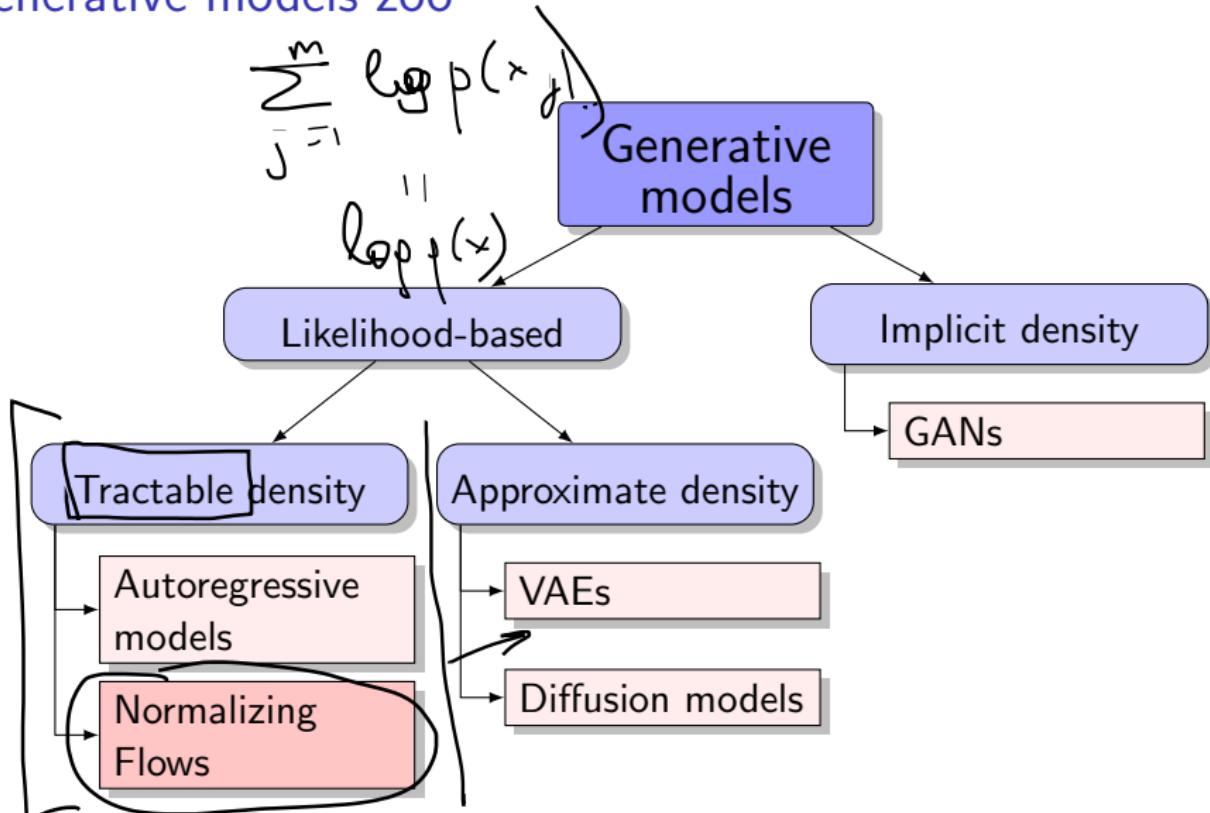
# Outline

1. Normalizing flows (NF)
2. NF examples
  - Linear normalizing flows
  - Gaussian autoregressive NF
  - Coupling layer (RealNVP)

# Outline

1. Normalizing flows (NF)
2. NF examples
  - Linear normalizing flows
  - Gaussian autoregressive NF
  - Coupling layer (RealNVP)

# Generative models zoo



# Normalizing flows prerequisites

$$\begin{cases} \cdot & \cdot \\ \cdot & \cdot \end{cases} = g$$

$$f(x)$$

$$g(z)$$

## Jacobian matrix

Let  $f : \mathbb{R}^m \rightarrow \mathbb{R}^m$  be a differentiable function.

$$z = f(x)$$

$$J = \frac{\partial z}{\partial x} = \begin{pmatrix} \frac{\partial z_1}{\partial x_1} & \cdots & \frac{\partial z_1}{\partial x_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial z_m}{\partial x_1} & \cdots & \frac{\partial z_m}{\partial x_m} \end{pmatrix} \in \mathbb{R}^{m \times m}$$

## Change of variable theorem (CoV)

Let  $x$  be a random variable with density function  $p(x)$  and

$f : \mathbb{R}^m \rightarrow \mathbb{R}^m$  is a differentiable, invertible function. If  $z = f(x)$ , then

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

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

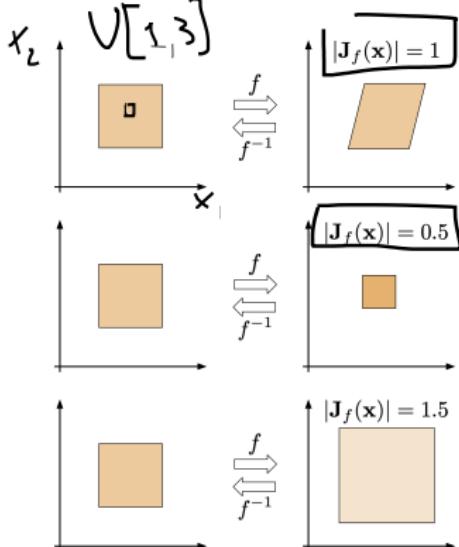
# Jacobian determinant

## Inverse function theorem

If function  $\mathbf{f}$  is invertible and Jacobian matrix is continuous and non-singular, then

$$\mathbf{J}_{\mathbf{f}^{-1}} = \mathbf{J}_{\mathbf{g}} = \mathbf{J}_{\mathbf{f}}^{-1}; \quad |\det(\mathbf{J}_{\mathbf{f}^{-1}})| = |\det(\mathbf{J}_{\mathbf{g}})| = \frac{1}{|\det(\mathbf{J}_{\mathbf{f}})|}.$$

- ▶  $\mathbf{x}$  and  $\mathbf{z}$  have the same dimensionality ( $\mathbb{R}^m$ ).
- ▶  $\mathbf{f}_{\theta}(\mathbf{x})$  could be parametric function. ( $NN$ )
- ▶ Determinant of Jacobian matrix  $\mathbf{J} = \frac{\partial \mathbf{f}_{\theta}(\mathbf{x})}{\partial \mathbf{x}}$  shows how the volume changes under the transformation.



# Fitting normalizing flows

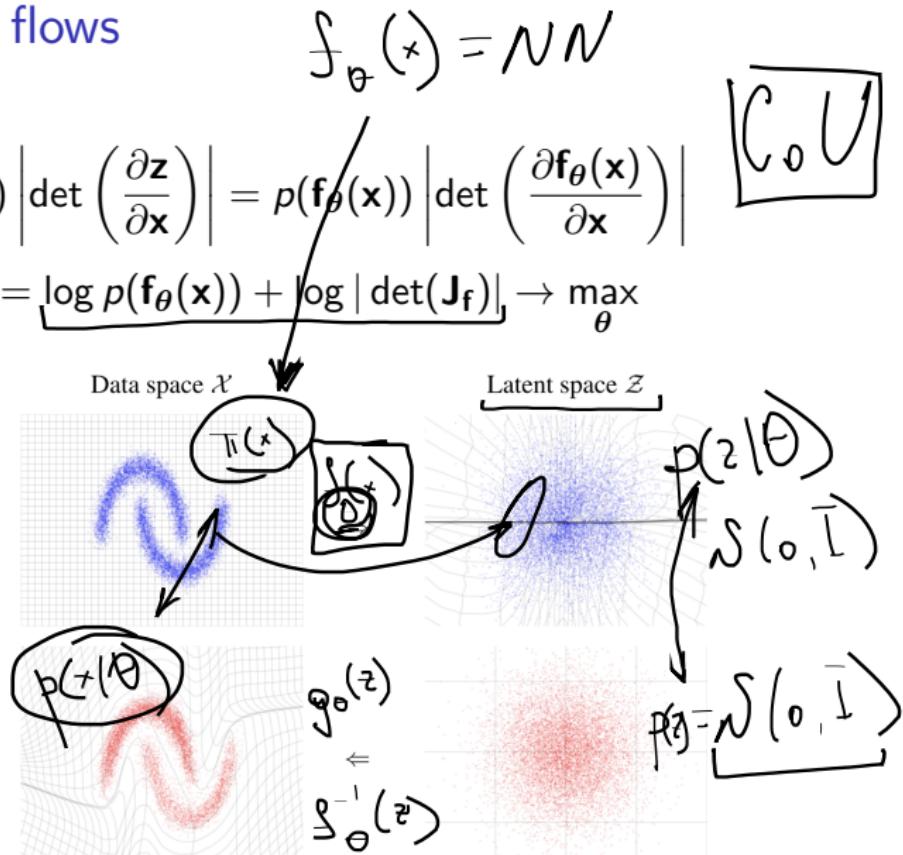
MLE problem

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

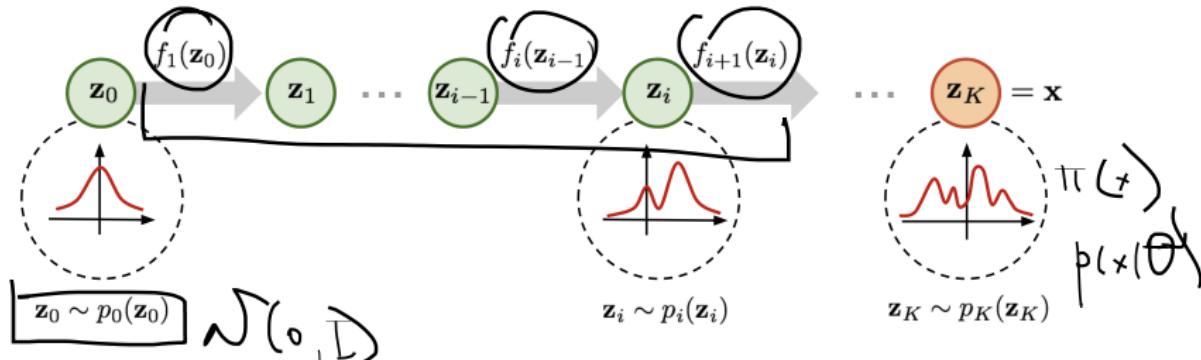
$$\log p(\mathbf{x}|\theta) = \underbrace{\log p(\mathbf{f}_\theta(\mathbf{x}))}_{\text{Data space } \mathcal{X}} + \underbrace{\log |\det(\mathbf{J}_f)|}_{\text{Latent space } \mathcal{Z}} \rightarrow \max_{\theta}$$

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

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



# Composition of normalizing flows



## Theorem

If  $\{\mathbf{f}_k\}_{k=1}^K$  satisfy conditions of the change of variable theorem, then  $\mathbf{z} = \mathbf{f}(\mathbf{x}) = \mathbf{f}_K \circ \dots \circ \mathbf{f}_1(\mathbf{x})$  also satisfies it.

$$\begin{aligned} p(\mathbf{x}) &= p(\mathbf{f}(\mathbf{x})) \left| \det \left( \frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} \right) \right| = p(\mathbf{f}(\mathbf{x})) \left| \det \left( \frac{\partial \mathbf{f}_K}{\partial \mathbf{f}_{K-1}} \cdots \frac{\partial \mathbf{f}_1}{\partial \mathbf{x}} \right) \right| = \\ &= p(\mathbf{f}(\mathbf{x})) \prod_{k=1}^K \left| \det \left( \frac{\partial \mathbf{f}_k}{\partial \mathbf{f}_{k-1}} \right) \right| = \boxed{p(\mathbf{f}(\mathbf{x}))} \prod_{k=1}^K |\det(\mathbf{J}_{f_k})| \end{aligned}$$

# Normalizing flows (NF)

Cat(i)

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

## Definition

Normalizing flow is a *differentiable, invertible* mapping from data  $\mathbf{x}$  to the noise  $\mathbf{z}$ .

- ▶ **Normalizing** means that NF takes samples from  $\pi(\mathbf{x})$  and normalizes them into samples from the base density  $p(\mathbf{z})$ .
- ▶ **Flow** refers to the trajectory followed by samples from  $p(\mathbf{z})$  as they are transformed by the sequence of transformations

$$\mathbf{z} = \mathbf{f}_K \circ \cdots \circ \mathbf{f}_1(\mathbf{x}); \quad \mathbf{x} = \mathbf{f}_1^{-1} \circ \cdots \circ \mathbf{f}_K^{-1}(\mathbf{z}) = \mathbf{g}_1 \circ \cdots \circ \mathbf{g}_K(\mathbf{z})$$

## Log likelihood

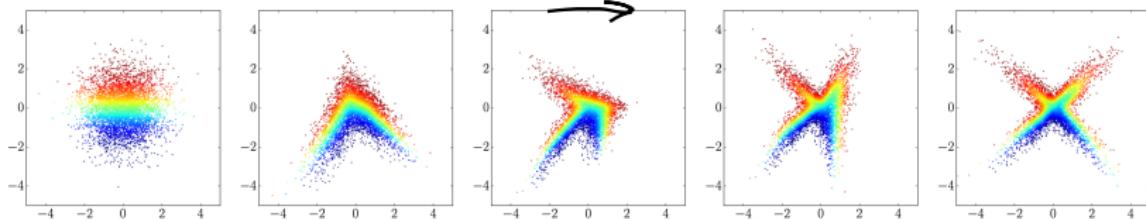
$$\log p(\underline{\mathbf{x}}|\theta) = \log p(\mathbf{f}_K \circ \cdots \circ \mathbf{f}_1(\mathbf{x})) + \sum_{k=1}^K \log |\det(\mathbf{J}_{\mathbf{f}_k})|,$$

where  $\mathbf{J}_{\mathbf{f}_k} = \frac{\partial \mathbf{f}_k}{\partial \mathbf{f}_{k-1}}$ .

**Note:** Here we consider only continuous random variables.

# Normalizing flows

Example of a 4-step NF



NF log likelihood

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

C.V  
 $\delta$   
 $x \in \mathbb{R}^d$

What is the complexity of the determinant computation?

What do we need?

- ▶ efficient computation of the Jacobian matrix  $\mathbf{J}_f = \frac{\partial \mathbf{f}_\theta(\mathbf{x})}{\partial \mathbf{x}}$ ;
- ▶ efficient inversion of  $\mathbf{f}_\theta(\mathbf{x})$ .

# Outline

## 1. Normalizing flows (NF)

## 2. NF examples

Linear normalizing flows

Gaussian autoregressive NF

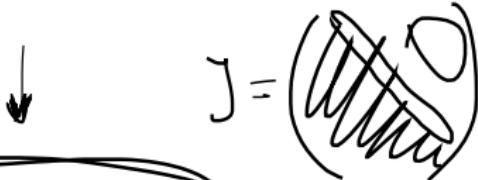
Coupling layer (RealNVP)

# Outline

1. Normalizing flows (NF)
2. NF examples
  - Linear normalizing flows
  - Gaussian autoregressive NF
  - Coupling layer (RealNVP)

## Jacobian structure

Normalizing flows log-likelihood



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

The main challenge is a determinant of the Jacobian matrix.

What is the  $\det(\mathbf{J})$  in the following cases?

Consider a linear layer  $\mathbf{z} = \mathbf{W}\mathbf{x}$ ,  $\mathbf{W} \in \mathbb{R}^{m \times m}$ .

1. Let  $\mathbf{z}$  be a permutation of  $\mathbf{x}$ .
2. Let  $z_j$  depend only on  $x_j$ .

$$\mathbf{z} = \begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_m \end{pmatrix} \quad O(m)$$

$$\log \left| \det \left( \frac{\partial \mathbf{f}_\theta(\mathbf{x})}{\partial \mathbf{x}} \right) \right| = \log \left| \prod_{j=1}^m \frac{\partial f_{j,\theta}(x_j)}{\partial x_j} \right| = \sum_{j=1}^m \log \left| \frac{\partial f_{j,\theta}(x_j)}{\partial x_j} \right|.$$

3. Let  $z_j$  depend only on  $\mathbf{x}_{1:j}$  (autoregressive dependency).

# Linear normalizing flows

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

In general, we need  $O(m^3)$  to invert matrix.

## Invertibility

- ▶ Diagonal matrix  $O(m)$ .
- ▶ Triangular matrix  $O(m^2)$ .
- ▶ It is impossible to parametrize all invertible matrices.

## Invertible 1x1 conv

$\mathbf{W} \in \mathbb{R}^{c \times c}$  – kernel of 1x1 convolution with  $c$  input and  $c$  output channels. The computational complexity of computing or differentiating  $\det(\mathbf{W})$  is  $O(c^3)$ . Cost to compute  $\det(\mathbf{W})$  is  $O(c^3)$ . It should be invertible.

# Linear normalizing flows

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

## Matrix decompositions

### ► LU-decomposition

$$\mathbf{W} = \mathbf{P}\mathbf{L}\mathbf{U},$$

where  $\mathbf{P}$  is a permutation matrix,  $\mathbf{L}$  is lower triangular with positive diagonal,  $\mathbf{U}$  is upper triangular with positive diagonal.

### ► QR-decomposition

$$\mathbf{W} = \boxed{\mathbf{Q}} \mathbf{R},$$

where  $\mathbf{Q}$  is an orthogonal matrix,  $\mathbf{R}$  is an upper triangular matrix with positive diagonal.

Decomposition should be done only once in the beginning. Next, we fit decomposed matrices ( $\mathbf{P}/\mathbf{L}/\mathbf{U}$  or  $\mathbf{Q}/\mathbf{R}$ ).

 Kingma D. P., Dhariwal P. Glow: Generative Flow with Invertible 1x1 Convolutions, 2018

Hoogeboom E., et al. Emerging convolutions for generative normalizing flows, 2019

# Outline

1. Normalizing flows (NF)
2. NF examples
  - Linear normalizing flows
  - Gaussian autoregressive NF
  - Coupling layer (RealNVP)

## Gaussian autoregressive model

C. V

Consider an autoregressive model

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

Sampling

$$\log p(\mathbf{x}|\theta) = \underbrace{\log p(\mathbf{x})}_{\text{Log p(x)}} + \underbrace{\log |\det(\mathbf{W})|}_{\text{Log |\det(W)|}}$$

$$x_j = \sigma_{j,\theta}(\mathbf{x}_{1:j-1}) z_j + \mu_{j,\theta}(\mathbf{x}_{1:j-1})$$

$$z_j \sim \mathcal{N}(0, 1)$$

Inverse transform

$$z_j = \frac{(x_j - \mu_{j,\theta}(\mathbf{x}_{1:j-1}))}{\sigma_{j,\theta}(\mathbf{x}_{1:j-1})}$$



- We have an invertible and differentiable transformation from  $p(\mathbf{z})$  to  $p(\mathbf{x}|\theta)$ .
- It is an autoregressive (AR) NF with the base distribution  $p(\mathbf{z}) = \mathcal{N}(0, \mathbf{I})$ !
- Jacobian of such transformation is triangular!

O(m)

## Gaussian autoregressive NF

$$\begin{cases} \mathbf{x} = \underline{\mathbf{g}_\theta(\mathbf{z})} \Rightarrow \underbrace{\mathbf{x}_j = \sigma_{j,\theta}(\mathbf{x}_{1:j-1}) \cdot z_j + \mu_{j,\theta}(\mathbf{x}_{1:j-1})}_{\text{Sampling}} \\ \mathbf{z} = \underline{\mathbf{f}_\theta(\mathbf{x})} \Rightarrow \underbrace{z_j = (\mathbf{x}_j - \mu_{j,\theta}(\mathbf{x}_{1:j-1})) \cdot \frac{1}{\sigma_{j,\theta}(\mathbf{x}_{1:j-1})}}_{\text{Training}} \end{cases}$$

Generation function  $\mathbf{g}_\theta(\mathbf{z})$  is sequential.

Inference function  $\underline{\mathbf{f}_\theta(\mathbf{x})}$  is **not sequential**.

### Forward KL for NF

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

- ▶ We need to be able to compute  $\mathbf{f}_\theta(\mathbf{x})$  and its Jacobian.
- ▶ We need to be able to compute the density  $p(\mathbf{z})$ .
- ▶ We don't need to think about computing the function  $\mathbf{g}_\theta(\mathbf{z}) = \mathbf{f}_\theta^{-1}(\mathbf{z})$  until we want to sample from the model.

# Gaussian autoregressive NF

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

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

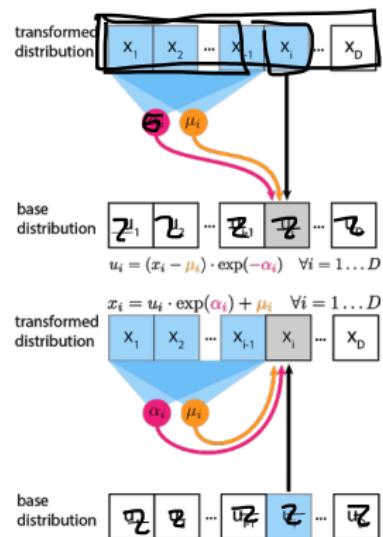
- ▶ Sampling is sequential, density estimation is parallel.
- ▶ Forward KL is a natural loss.

Forward transform:  $\mathbf{f}_\theta(\mathbf{x})$     *train*

$$z_j = (x_j - \mu_{j,\theta}(\mathbf{x}_{1:j-1})) \cdot \frac{1}{\sigma_{j,\theta}(\mathbf{x}_{1:j-1})}$$

Inverse transform:  $\mathbf{g}_\theta(\mathbf{z})$

$$x_j = \sigma_{j,\theta}(\mathbf{x}_{1:j-1}) \cdot z_j + \mu_{j,\theta}(\mathbf{x}_{1:j-1})$$



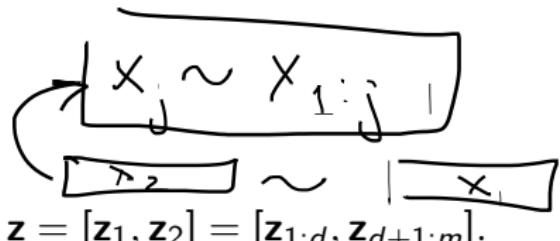
# Outline

1. Normalizing flows (NF)
2. NF examples
  - Linear normalizing flows
  - Gaussian autoregressive NF
  - Coupling layer (RealNVP)

# RealNVP

Let split  $\mathbf{x}$  and  $\mathbf{z}$  in two parts:

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

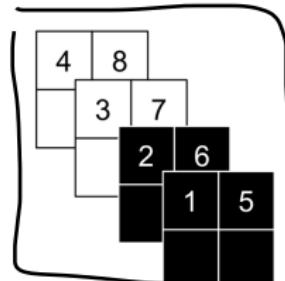


## Coupling layer

$$\begin{cases} \mathbf{x}_1 = \mathbf{z}_1; \\ \mathbf{x}_2 = \mathbf{z}_2 \odot \sigma_\theta(\mathbf{z}_1) + \mu_\theta(\mathbf{z}_1). \end{cases}$$

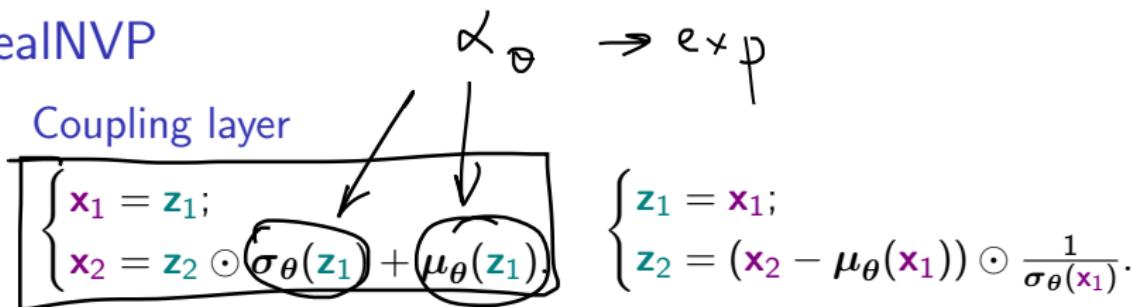
$$\begin{cases} \mathbf{z}_1 = \mathbf{x}_1; \\ \mathbf{z}_2 = (\mathbf{x}_2 - \mu_\theta(\mathbf{x}_1)) \odot \frac{1}{\sigma_\theta(\mathbf{x}_1)}. \end{cases}$$

## Image partitioning



- Checkerboard ordering uses masking.
- Channelwise ordering uses splitting.

## RealNVP



Estimating the density takes 1 pass, sampling takes 1 pass!

## Jacobian

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

$\mathcal{O}(m)$

## Gaussian AR NF

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

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

How to get RealNVP coupling layer from gaussian AR NF?

## Glow samples

Glow model: [coupling layer] + [linear flows (1x1 convs)]

Flow match  
ng



Kingma D. P., Dhariwal P. Glow: Generative Flow with Invertible 1x1 Convolutions,

2018

## Summary

- ▶ Change of variable theorem allows to get the density function of the random variable under the invertible transformation.
- ▶ Normalizing flows transform a simple base distribution to a complex one via a sequence of invertible transformations with tractable Jacobian.
- ▶ Normalizing flows have a tractable likelihood that is given by the change of variable theorem.
- ▶ Linear NF try to parametrize set of invertible matrices via matrix decompositions.
- ▶ Gaussian autoregressive NF is an autoregressive model with triangular Jacobian.
- ▶ The RealNVP coupling layer is an effective type of NF (special case of AR NF) that has fast inference and generation modes.