# Deep Generative Models

## Lecture 2

Roman Isachenko

Moscow Institute of Physics and Technology
Yandex School of Data Analysis

2025, Autumn

# Recap of Previous Lecture

We're given i.i.d. samples $\{\mathbf{x}_i\}_{i=1}^{n} \subset \mathbb{R}^m$ drawn from some unknown distribution $\pi(\mathbf{x})$.

## Objective

Our goal is to learn the distribution $\pi(\mathbf{x})$ so that we can:

- Evaluate $\pi(\mathbf{x})$ for new samples;
- Sample from $\pi(\mathbf{x})$ (i.e., generate novel samples $\mathbf{x} \sim \pi(\mathbf{x})$).

Rather than considering all possible probability distributions, we approximate $\pi(\mathbf{x})$ by a parameterized family $p(\mathbf{x}|\boldsymbol{\theta}) \approx \pi(\mathbf{x})$.

## Divergence Minimization Task

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

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

# Recap of Previous Lecture

### Forward KL Divergence

$$\mathrm{KL}(\pi\|p) = \int \pi(\mathbf{x}) \log \frac{\pi(\mathbf{x})}{p(\mathbf{x}|\boldsymbol{\theta})} \, d\mathbf{x} \to \min_{\boldsymbol{\theta}}$$

### Reverse KL Divergence

$$\mathrm{KL}(p\|\pi) = \int p(\mathbf{x}|\boldsymbol{\theta}) \log \frac{p(\mathbf{x}|\boldsymbol{\theta})}{\pi(\mathbf{x})} \, d\mathbf{x} \to \min_{\boldsymbol{\theta}}$$

### Maximum Likelihood Estimation (MLE)

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

Maximum likelihood estimation is equivalent to minimizing the Monte Carlo estimate of the forward KL divergence.

# Recap of Previous Lecture

### Likelihood as Product of Conditionals

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

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

### MLE for Autoregressive Models

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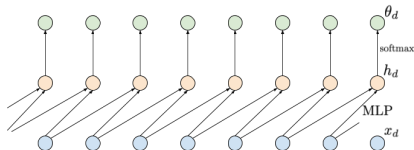
### Sampling

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

The generated sample is $\hat{\mathbf{x}} = (\hat{x}_1, \hat{x}_2, \ldots, \hat{x}_m)$.

# Recap of Previous Lecture

## Autoregressive MLP
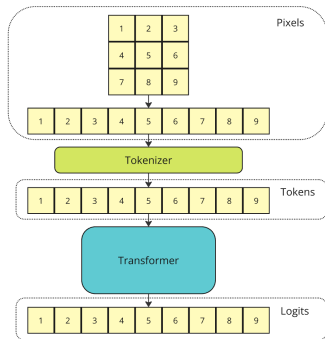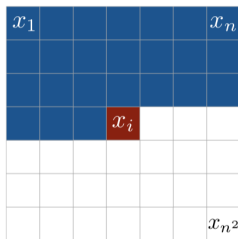


## Autoregressive Transformer



Image credit: https://jmtomczak.github.io/blog/2/2_ARM.html
Chen M. et al. Generative Pretraining from Pixels, 2020

# Outline

# Outline

# Generative Models Zoo

# Normalizing Flows: Prerequisites

### Jacobian Matrix

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

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

### Change of Variables Theorem (CoV)

Let $\mathbf{x}$ be a random variable with density $p(\mathbf{x})$ and $\mathbf{f} : \mathbb{R}^m \to \mathbb{R}^m$ a differentiable, **invertible** mapping. If $\mathbf{z} = \mathbf{f}(\mathbf{x})$ and $\mathbf{x} = \mathbf{f}^{-1}(\mathbf{z}) = \mathbf{g}(\mathbf{z})$, then

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

$$p(\mathbf{z}) = p(\mathbf{x})|\det(\mathbf{J_g})| = p(\mathbf{x})\left|\det\left(\frac{\partial \mathbf{x}}{\partial \mathbf{z}}\right)\right| = p(\mathbf{g}(\mathbf{z}))\left|\det\left(\frac{\partial \mathbf{g}(\mathbf{z})}{\partial \mathbf{z}}\right)\right|$$
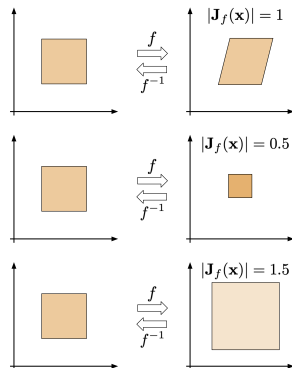
# Jacobian Determinant

### Inverse Function Theorem

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

$$\mathbf{J_{f^{-1}}} = \mathbf{J_g} = \mathbf{J_f^{-1}}; \quad |\det(\mathbf{J_{f^{-1}}})| = |\det(\mathbf{J_g})| = \frac{1}{|\det(\mathbf{J_f})|}$$

- $\mathbf{x}$ and $\mathbf{z}$ reside in the same space $(\mathbb{R}^m)$.

- $\mathbf{f_\theta}(\mathbf{x})$ is a parameterized transformation.

- The determinant of the Jacobian $\mathbf{J} = \frac{\partial \mathbf{f_\theta}(\mathbf{x})}{\partial \mathbf{x}}$ quantifies how the volume is changed by the transformation.
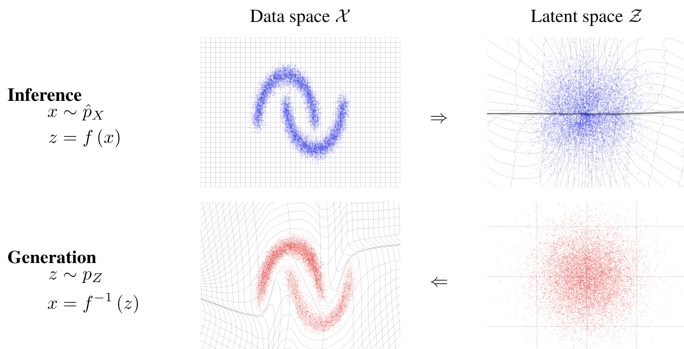


https://jmtomczak.github.io/blog/3/3_flows.html

# Fitting Normalizing Flows

## MLE Problem

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

$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \log p(\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})) + \log|\det(\mathbf{J_f})| \to \max_{\boldsymbol{\theta}}$$



Dinh L., Sohl-Dickstein J., Bengio S. Density Estimation Using Real NVP, 2016

# Composition of Normalizing Flows



$$z_0 \xrightarrow{f_1(z_0)} z_1 \cdots z_{i-1} \xrightarrow{f_i(z_{i-1})} z_i \xrightarrow{f_{i+1}(z_i)} \cdots z_K = x$$

$$z_0 \sim p_0(z_0) \qquad z_i \sim p_i(z_i) \qquad z_K \sim p_K(z_K)$$

### Theorem

If every $\{f_k\}_{k=1}^K$ satisfies the conditions of the change-of-variables theorem, then the composition $f(x) = f_K \circ \ldots \circ f_1(x)$ also satisfies them.

$$p(x) = p(f(x)) \left| \det\left( \frac{\partial f(x)}{\partial x} \right) \right| = p(f(x)) \left| \det\left( \frac{\partial f_K}{\partial f_{K-1}} \ldots \frac{\partial f_1}{\partial x} \right) \right| =$$

$$= p(f(x)) \prod_{k=1}^K \left| \det\left( \frac{\partial f_k}{\partial f_{k-1}} \right) \right| = p(f(x)) \prod_{k=1}^K |\det(J_{f_k})|$$

# Normalizing Flows (NF)

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

### Definition

A normalizing flow is a *differentiable, invertible* mapping that transforms data $\mathbf{x}$ to latent noise $\mathbf{z}$.

▶ **Normalizing** refers to mapping samples from $\pi(\mathbf{x})$ to a base distribution $p(\mathbf{z})$.

▶ **Flow** describes the sequence of transformations that maps samples from $p(\mathbf{z})$ to the target, more complex distribution.

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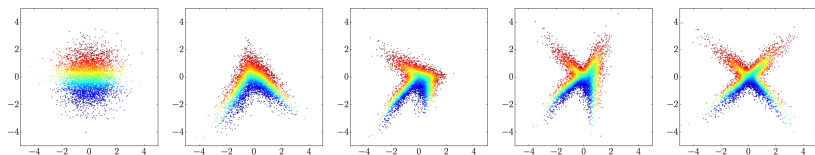
### Log-Likelihood

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

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: 4-Step NF



## NF Log-Likelihood

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

What's the computational complexity of evaluating this determinant?

## Requirements

▶ Efficient computation of the Jacobian $\mathbf{J_f} = \frac{\partial \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})}{\partial \mathbf{x}}$

▶ Efficient inversion of the transformation $\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})$

*Papamakarios G. et al. Normalizing Flows for Probabilistic Modeling and Inference, 2019*

# Outline

# Outline

# Jacobian Structure

### Normalizing Flows Log-Likelihood

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

The principal computational challenge is evaluating the Jacobian determinant.

### What is det($\mathbf{J}$) in These Cases?

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

1. $\mathbf{z}$ is a permutation of $\mathbf{x}$.
2. $z_j$ depends only on $x_j$.

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

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

# Linear Normalizing Flows

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

In general, matrix inversion has computational complexity $O(m^3)$.

Invertibility

- ▶ Diagonal matrix: $O(m)$.
- ▶ Triangular matrix: $O(m^2)$.
- ▶ Directly parameterizing the full group of invertible matrices is infeasible.

### Invertible $1 \times 1$ Convolution

$\mathbf{W} \in \mathbb{R}^{c \times c}$ acts as the kernel of a $1 \times 1$ convolution with $c$ input and $c$ output channels. Calculating or differentiating $\det(\mathbf{W})$ incurs a cost of $O(c^3)$. It is critical that $\mathbf{W}$ is invertible.

---

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

# Linear Normalizing Flows

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

## Matrix Decompositions

▶ **LU Decomposition:**

$$\mathbf{W} = \mathbf{PLU},$$

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

▶ **QR Decomposition:**

$$\mathbf{W} = \mathbf{QR},$$

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

Decomposition is performed only at initialization; the decomposed matrices ($\mathbf{P}, \mathbf{L}, \mathbf{U}$ or $\mathbf{Q}, \mathbf{R}$) are optimized during training.

*Kingma D. P., et al. Glow: Generative Flow with Invertible 1x1 Convolutions, 2018*
*Hoogeboom E., et al. Emerging Convolutions for Generative Normalizing Flows, 2019*

# Outline

# Gaussian Autoregressive Model

Consider the autoregressive model:

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

## Sampling

$$x_j = \sigma_{j,\boldsymbol{\theta}}(\mathbf{x}_{1:j-1}) \cdot z_j + \mu_{j,\boldsymbol{\theta}}(\mathbf{x}_{1:j-1}), \quad z_j \sim \mathcal{N}(0,1)$$

## Inverse Transformation

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

▶ This gives an **invertible** and **differentiable** transformation from $p(\mathbf{z})$ to $p(\mathbf{x}|\boldsymbol{\theta})$.

▶ This model is called an autoregressive (AR) NF with base distribution $p(\mathbf{z}) = \mathcal{N}(0, \mathbf{I})$.

▶ The Jacobian matrix of this transformation is triangular.

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

# Gaussian Autoregressive NF

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

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

To generate samples, apply $\mathbf{g}_\theta(\mathbf{z})$ sequentially;
inference via $\mathbf{f}_\theta(\mathbf{x})$ is parallelizable.

## Forward KL for NFs

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

- ▶ Computing $\mathbf{f}_\theta(\mathbf{x})$ and its Jacobian is necessary.
- ▶ One must be able to evaluate the density $p(\mathbf{z})$.
- ▶ The inverse $\mathbf{g}_\theta(\mathbf{z}) = \mathbf{f}_\theta^{-1}(\mathbf{z})$ is only needed for sampling.

---

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

# Gaussian Autoregressive NF

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

$$\mathbf{z} = \mathbf{f}_\theta(\mathbf{x}) \quad \Rightarrow \quad z_j = \frac{x_j - \mu_{j,\theta}(\mathbf{x}_{1:j-1})}{\sigma_{j,\theta}(\mathbf{x}_{1:j-1})}$$
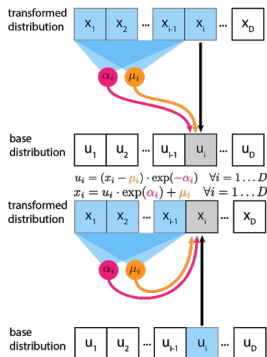
▶ Sampling must be done sequentially, but density estimation can be parallelized.

▶ The forward KL divergence is a natural objective for training.

Forward Transformation: $\mathbf{f}_\theta(\mathbf{x})$

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

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

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



*Image credit: https://blog.evjang.com/2018/01/nf2.html*

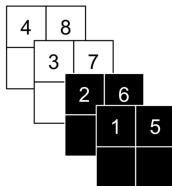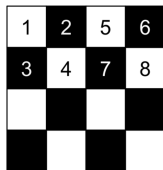# Outline

# RealNVP

Split **x** and **z** into 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

$$\begin{cases} \mathbf{x}_1 = \mathbf{z}_1 \\ \mathbf{x}_2 = \mathbf{z}_2 \odot \boldsymbol{\sigma}_{\boldsymbol{\theta}}(\mathbf{z}_1) + \boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{z}_1) \end{cases} \qquad \begin{cases} \mathbf{z}_1 = \mathbf{x}_1 \\ \mathbf{z}_2 = (\mathbf{x}_2 - \boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{x}_1)) \odot \frac{1}{\boldsymbol{\sigma}_{\boldsymbol{\theta}}(\mathbf{x}_1)} \end{cases}$$

## Image Partitioning



▶ Checkerboard ordering corresponds to masking.

▶ Channelwise ordering relies on splitting.

*Dinh L., Sohl-Dickstein J., Bengio S. Density Estimation Using Real NVP, 2016*

# RealNVP

### Coupling Layer

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

In both training and sampling, only a single forward pass is needed!

### Jacobian

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

### 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 can the RealNVP layer be derived as a special instance of the
<u>Gaussian autoregressive</u> NF?

*Dinh L., Sohl-Dickstein J., Bengio S. Density Estimation Using Real NVP, 2016*

# Glow: Coupling Layers + Linear Flows ($1 \times 1$ Convolutions)



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

# Summary

▶ The change-of-variables theorem provides a method for computing a random variable's density under an invertible transformation.

▶ Normalizing flows transform a simple base distribution into a complex one via a sequence of invertible mappings, each with efficient Jacobian determinants.

▶ This enables exact likelihood computation, thanks to the change-of-variables formula.

▶ Linear NFs capture invertible matrices by using matrix decompositions.

▶ Gaussian autoregressive NFs are AR models with triangular Jacobians.

▶ The RealNVP coupling layer provides an efficient normalizing flow (a special case of AR NF), supporting fast inference and sampling.