# VI & BNN (2)

# Flow-based Models

Keywords : Normalizing Flow, Masked Autoregressive Flow,

Inverse Autoregressive Flow, Glow, RealNVP, NICE

Seunghan Lee

2021.03.15.Mon

# Contents

1. What is Normalizing Flow?
   1. Generative Modeling
   2. Change of Variables
   3. What is Normalizing Flow?

2. Models using Normalizing Flows
   1. Real NVP (Real-valued Non-Volume Preserving)
   2. NICE (Non-linear Independent Component Estimation)
   3. Glow (Generative Flow)

3. Models using Autoregressive Flows
   1. MADE (Masked Autoencoder for Density Estimation)
   2. MAF (Masked Autoregressive Flows)
   3. IAF (Inverse Autoregressive Flows)

# Paper List / References

**Papers**

- Non-linear Independent Components Estimation (NICE) ( Dinh et al., 2014 )

- MADE : Masked Autoencoder for Distribution Estimation ( Germain et al., 2015)

- Variational Inference with Normalizing Flows ( Rezende and Mohamed, 2016 )

- Improved Variational Inference with Inverse Autoregressive Flow ( Kingma et al., 2016)

- Density Estimation using Real NVP ( Dinh et al., 2017 )

- Masked Autoregressive Flow for Density Estimation ( Papamakarios et al., 2017)

- Glow : Generative Flow with Invertible 1x1 Convolutions ( Kingma and Dhariwal, 2018)

**Blog**

- **https://seunghan96.github.io/**

- https://lilianweng.github.io/lil-log/2018/10/13/flow-based-deep-generative-models.html#glow

- https://sites.google.com/view/berkeley-cs294-158-sp20/home

# 1. What is Normalizing Flow?

# 1. What is Normalizing Flow?

## 1-1. Generative Modeling

Generative Model vs Discriminative Model

**Generative Model** : captures the joint probability $P(X, Y)$

( or just $P(X)$ if Y does not exists )

**Discriminative Model** : captures the conditional probability $P(Y \mid X)$

# 1. What is Normalizing Flow?

## 1-1. Generative Modeling

Generative Model vs Discriminative Model

**Generative Model** : captures the joint probability $P(X, Y)$

( or just $P(X)$ if Y does not exists )

**Discriminative Model** : captures the conditional probability $P(Y \mid X)$

# 1. What is Normalizing Flow?

## 1-1. Generative Modeling

Generative Model vs Discriminative Model

**Generative Model** : captures the joint probability $P(X, Y)$

( or just $P(X)$ if Y does not exists )

**Discriminative Model** : captures the conditional probability $P(Y \mid X)$

What do we want to do from this Generative Model?

1) Sampling

2) Density Evaluation

# 1. What is Normalizing Flow?

## 1-1. Generative Modeling

**Generative Model** : models that can generate samples (data)!

ex) GAN, VAE

- both shows an impressive result on hard tasks, such as images!

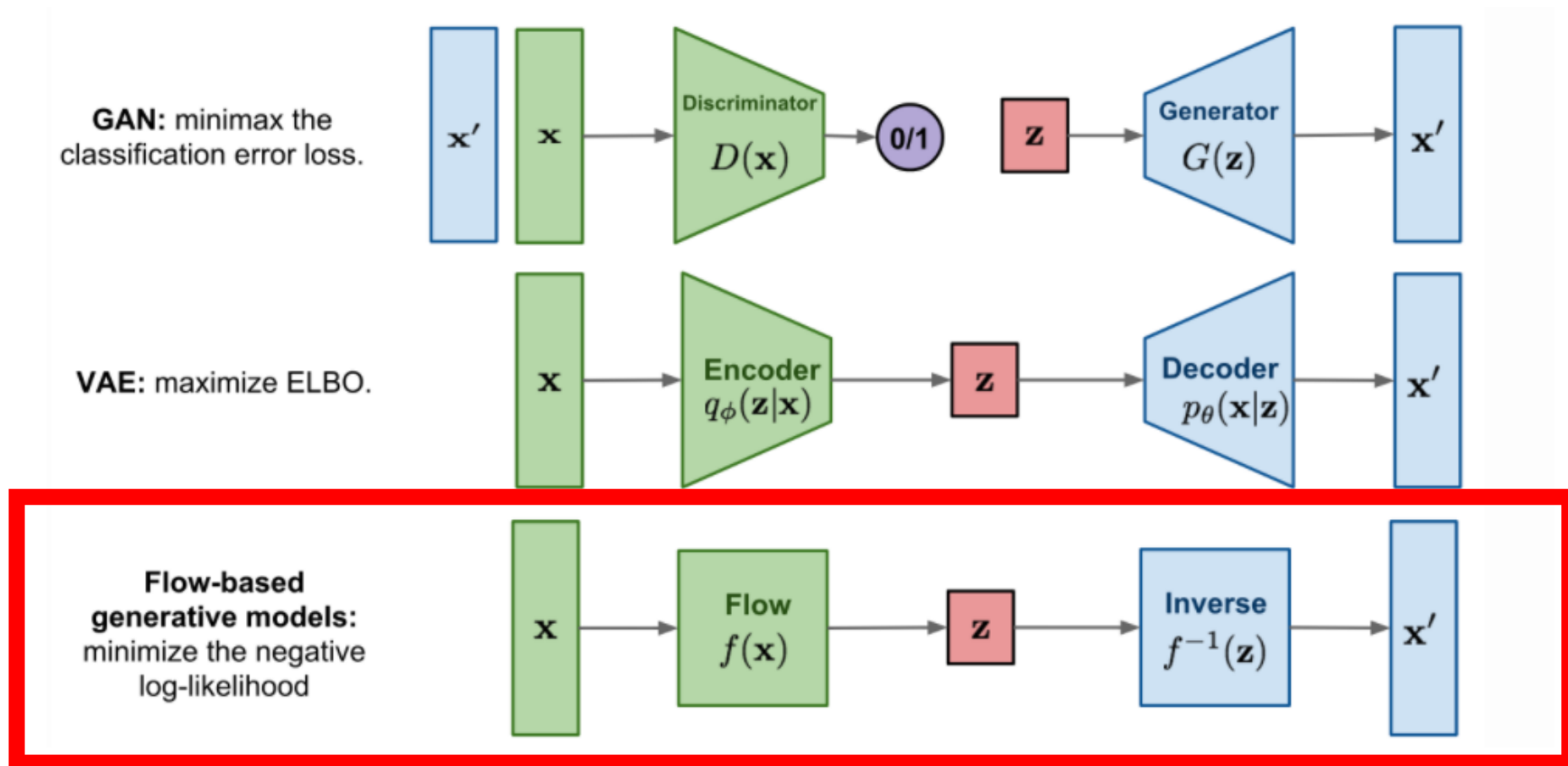- But cannot evaluate exact probability density of new points!


Using **Normalizing Flow**...

we can do both **"(1) sampling" & "(2) density evaluation"** efficiently & exacty!


So, what is Normalizing Flow?
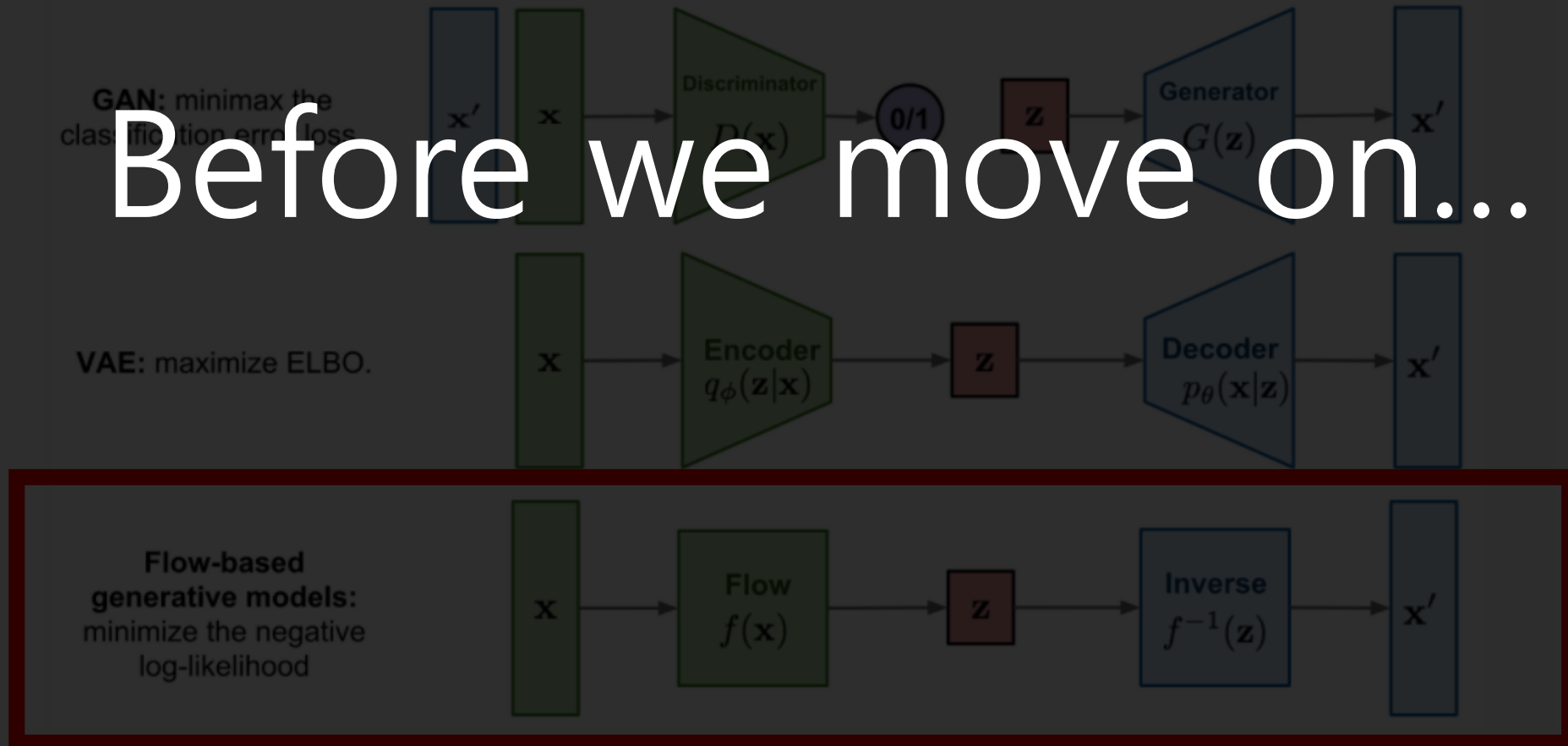
# 1. What is Normalizing Flow?

## 1-1. Generative Modeling

**Generative Model** : models that can generate samples (data)!

ex) GAN, VAE

- both shows ░(1) Generate new image! ░ hard░

- But cannot ░░░░ty de░

░(2) Find the density of

a given image!░

Using **Normalizing Flow**...

we can do both "**(1) sampling**" & "**(2) density evaluation"** efficiently & exacty!

So, what is Normalizing Flow?

# 1. What is Normalizing Flow?

## 1-1. Generative Modeling

## 1-1. Generative Modeling



Before we move on...

# 1. What is Normalizing Flow?

## 1-2. Change of Variables

Univariate Case

$$\int p(x)dx = \int \pi(z)dz = 1 \text{ ; Definition of probability distribution.}$$

$$p(x) = \pi(z)\left|\frac{dz}{dx}\right| = \pi(f^{-1}(x))\left|\frac{df^{-1}}{dx}\right| = \pi(f^{-1}(x))|(f^{-1})'(x)|$$

Multivariate Case

$$\mathbf{z} \sim \pi(\mathbf{z}), \mathbf{x} = f(\mathbf{z}), \mathbf{z} = f^{-1}(\mathbf{x})$$

$$p(\mathbf{x}) = \pi(\mathbf{z})\left|\det \frac{d\mathbf{z}}{d\mathbf{x}}\right| = \pi(f^{-1}(\mathbf{x}))\left|\det \frac{df^{-1}}{d\mathbf{x}}\right|$$

# 1. What is Normalizing Flow?

## 1-2. Change of Variables

Univariate Case

$$\int p(x)dx = \int \pi(z)dz = 1 \text{ ; Definition of probability distribution.}$$

$$p(x) = \pi(z)\left|\frac{dz}{dx}\right| = \pi(f^{-1}(x))\left|\frac{df^{-1}}{dx}\right| = \pi(f^{-1}(x))|(f^{-1})'(x)|$$

Multivariate Case

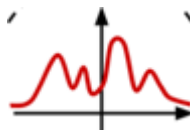$$\mathbf{z} \sim \pi(\mathbf{z}), \mathbf{x} = f(\mathbf{z}), \mathbf{z} = f^{-1}(\mathbf{x})$$

$$p(\mathbf{x}) = \pi(\mathbf{z})\left|\det\frac{d\mathbf{z}}{d\mathbf{x}}\right| = \pi(f^{-1}(\mathbf{x}))\left|\det\frac{df^{-1}}{d\mathbf{x}}\right|$$

Key point :

Need to calculate the **Jacobian term!**

# 1. What is Normalizing Flow?

## 1-3. What is Normalizing Flow?

Goal : We want to find out the unknown, complex pdf $p(x)$
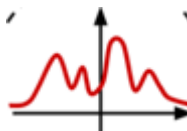
# 1. What is Normalizing Flow?

## 1-3. What is Normalizing Flow?

Goal : We want to find out the unknown, complex pdf $p(x)$ 

1)  Want to **SAMPLE IMAGE** from this distribution!

2)  Want to **EVALUATE the DENSITY** of the image!

# 1. What is Normalizing Flow?
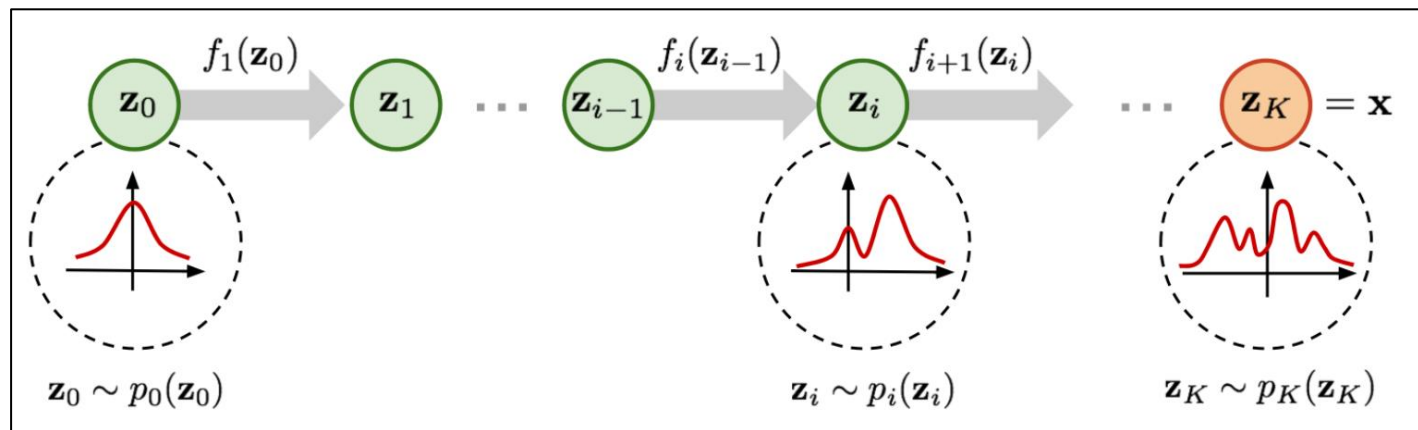
## 1-3. What is Normalizing Flow?

Goal : We want to find out the unknown, complex pdf $p(x)$ 

Intuitive idea of NF :

(1) Make a simple distribution!

(2) Change x N until it becomes the distribution that we want! ( sequentially )
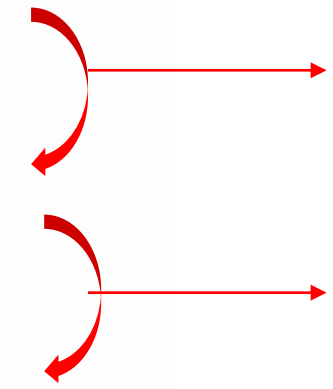
( from $p_0(\mathbf{z}_0)$ to $p_K(\mathbf{z}_K)$ )

# 1. What is Normalizing Flow?

## 1-3. What is Normalizing Flow?

By the <span style="color:red">change of variables</span> formula..

$$\mathbf{z}_{i-1} \sim p_{i-1}(\mathbf{z}_{i-1})$$

$$\mathbf{z}_i = f_i(\mathbf{z}_{i-1}), \text{ thus } \mathbf{z}_{i-1} = f_i^{-1}(\mathbf{z}_i)$$

$$p_i(\mathbf{z}_i) = p_{i-1}(f_i^{-1}(\mathbf{z}_i)) \left| \det \frac{df_i^{-1}}{d\mathbf{z}_i} \right|$$

$$= p_{i-1}(\mathbf{z}_{i-1}) \left| \det \left( \frac{df_i}{d\mathbf{z}_{i-1}} \right)^{-1} \right|$$

$$= p_{i-1}(\mathbf{z}_{i-1}) \left| \det \frac{df_i}{d\mathbf{z}_{i-1}} \right|^{-1}$$

$$\frac{df^{-1}(y)}{dy} = \frac{dx}{dy} = \left(\frac{dy}{dx}\right)^{-1} = \left(\frac{df(x)}{dx}\right)^{-1}$$

$$\det(M^{-1}) = (\det(M))^{-1}$$

# 1. What is Normalizing Flow?

## 1-3. What is Normalizing Flow?

$$\log p_i(\mathbf{z}_i) = \log p_{i-1}(\mathbf{z}_{i-1}) - \log\left|\det \frac{df_i}{d\mathbf{z}_{i-1}}\right|$$
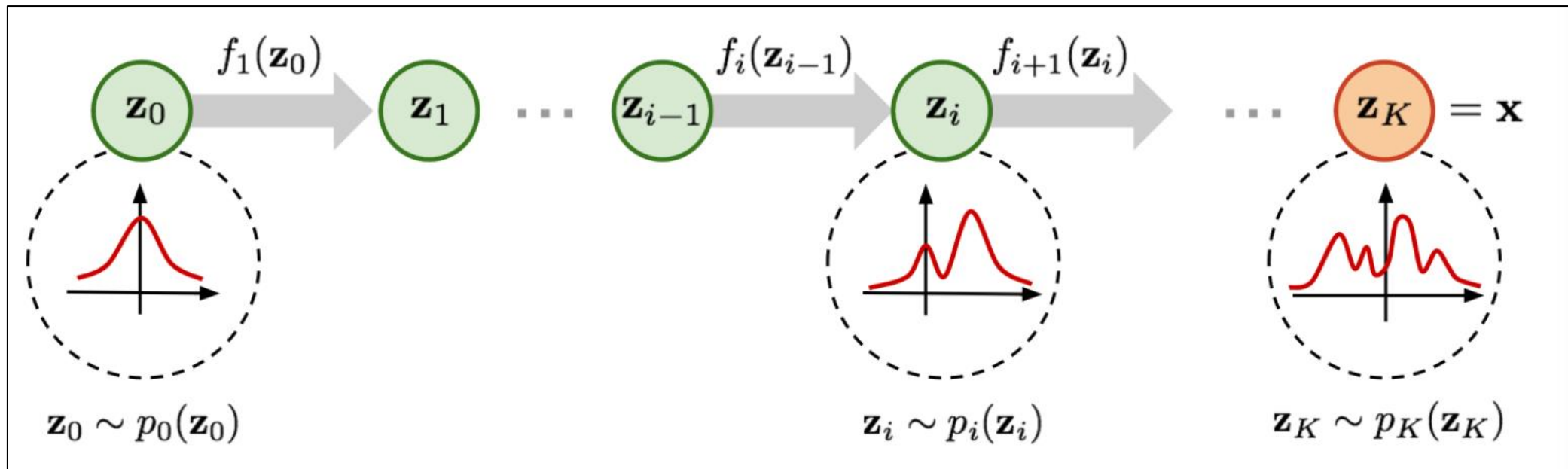
Sequentially apply the multiple functions to our base distribution!

$$\mathbf{x} = \mathbf{z}_K = f_K \circ f_{K-1} \circ \cdots \circ f_1(\mathbf{z}_0)$$

$$
\begin{aligned}
\log p(\mathbf{x}) = \log \pi_K(\mathbf{z}_K) &= \log \pi_{K-1}(\mathbf{z}_{K-1}) - \log\left|\det \frac{df_K}{d\mathbf{z}_{K-1}}\right| \\
&= \log \pi_{K-2}(\mathbf{z}_{K-2}) - \log\left|\det \frac{df_{K-1}}{d\mathbf{z}_{K-2}}\right| - \log\left|\det \frac{df_K}{d\mathbf{z}_{K-1}}\right| \\
&= \ldots \\
&= \log \pi_0(\mathbf{z}_0) - \sum_{i=1}^{K} \log\left|\det \frac{df_i}{d\mathbf{z}_{i-1}}\right|
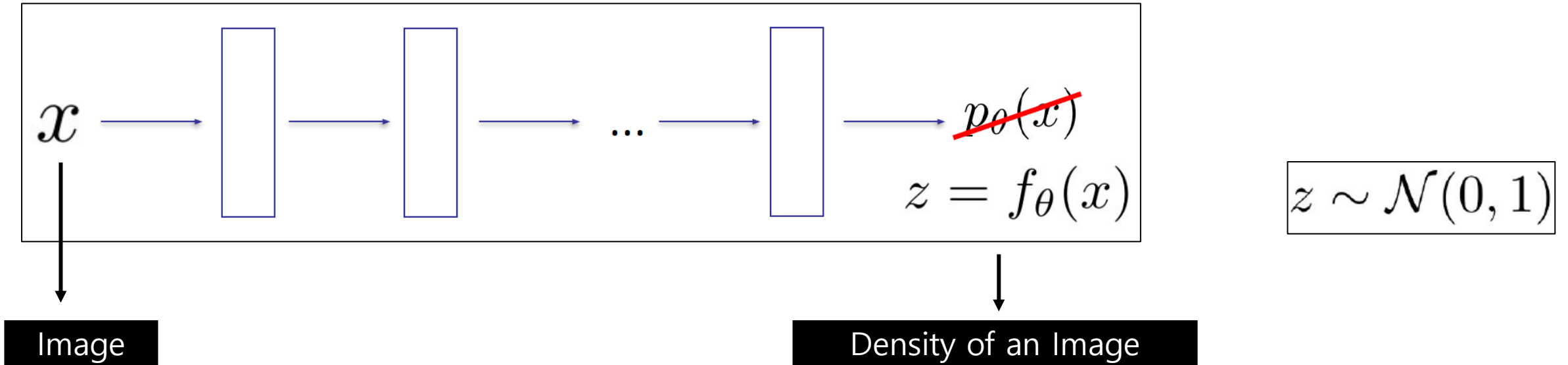\end{aligned}
$$

## 1-3. What is Normalizing Flow?
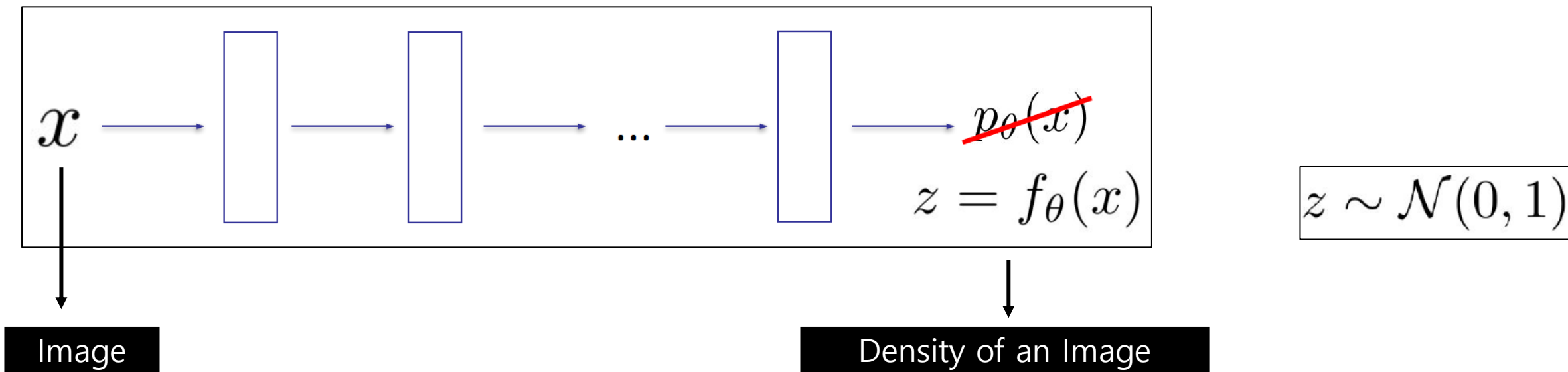
# 1. What is Normalizing Flow?

## 1-4. How to train NF?



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

Image

Density of an Image

# 1. What is Normalizing Flow?

## 1-4. How to train NF?



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

Image

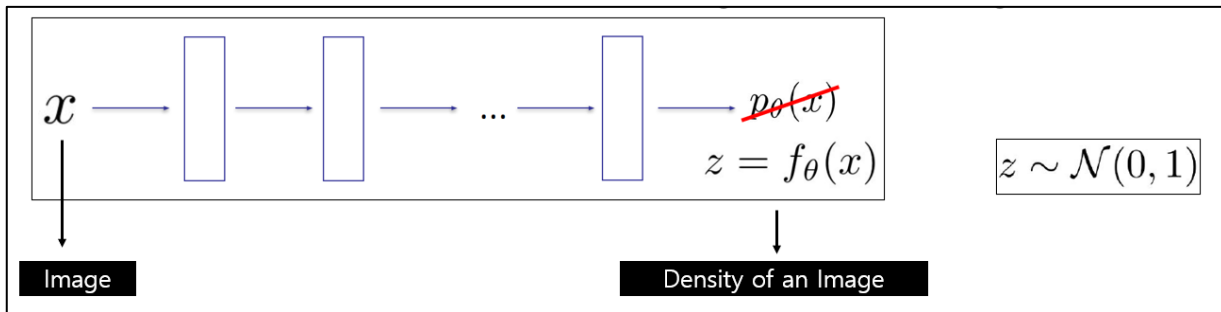Density of an Image

How do we train this model?

# 1. What is Normalizing Flow?

## 1-4. How to train NF?



$$z^{(i)} = f_\theta(x^{(i)})$$

$$p_\theta(x^{(i)}) = p_Z(z^{(i)}) \left| \frac{\partial z}{\partial x}(x^{(i)}) \right|$$

$$= p_Z(f_\theta(x^{(i)})) \left| \frac{\partial f_\theta}{\partial x}(x^{(i)}) \right|$$

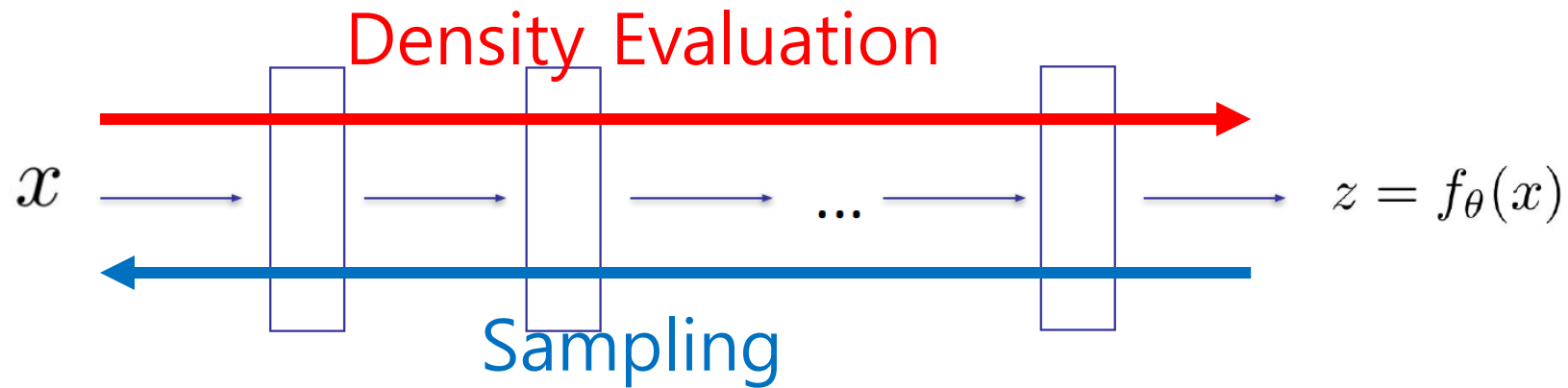$$\max_\theta \sum_i \log p_\theta(x^{(i)}) \implies \max_\theta \sum_i \log p_Z(f_\theta(x^{(i)})) + \log \left| \frac{\partial f_\theta}{\partial x}(x^{(i)}) \right|$$

**Train this using SGD!**

# 1. What is Normalizing Flow?

## 1-4. How to train NF?
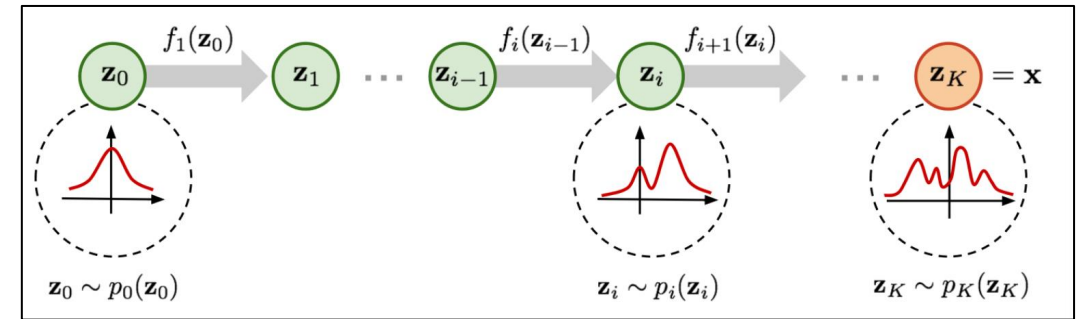
# 1. What is Normalizing Flow?

## 1-4. How to train NF?

So, what should we choose for our function $f_i$ ?



1. Easily Invertible

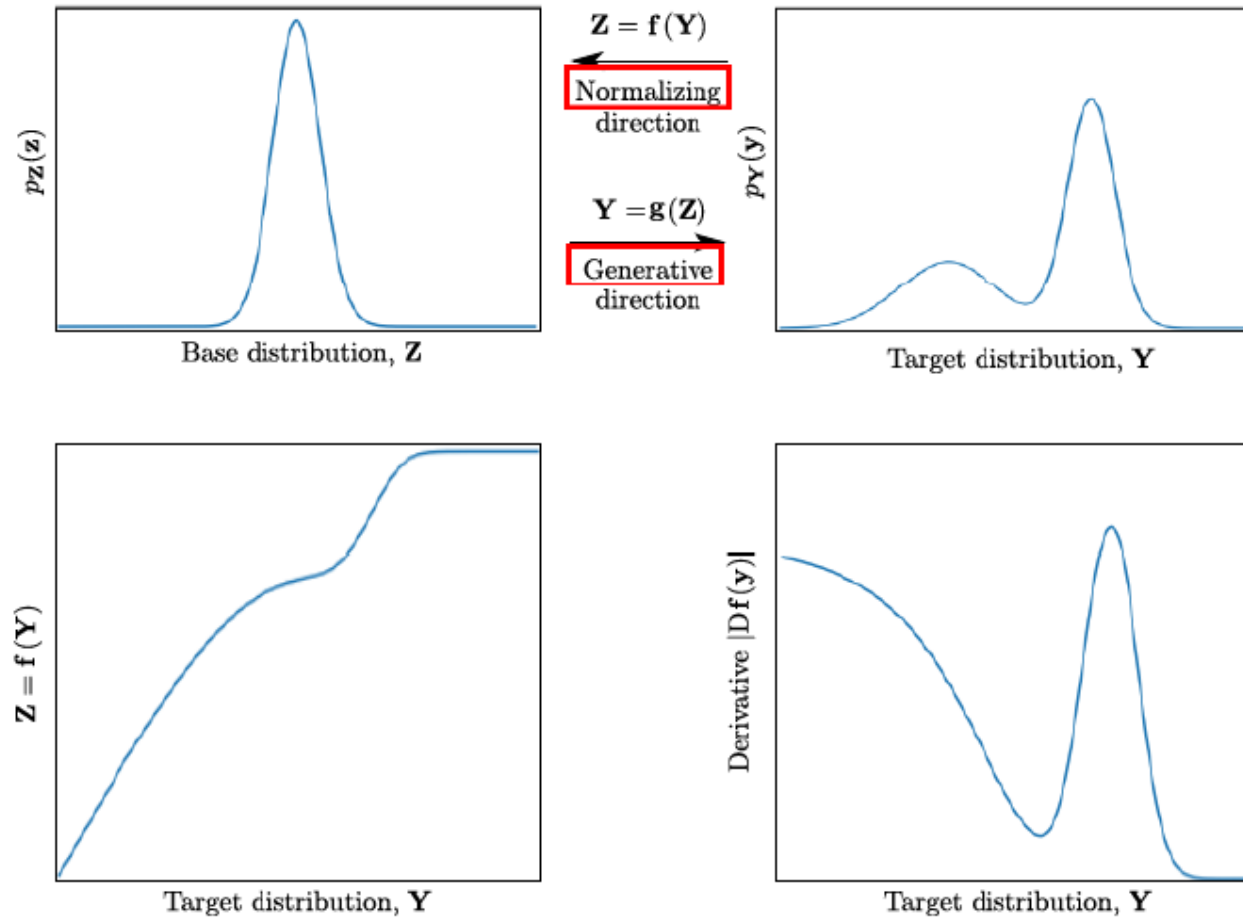2. Determinant of Jacobian to be easily computed

( Computation time for finding the Jacobian for matrix of D-dim : O(D^3) )

Lots of algorithms have been proposed that meets those two goals!

# 2. Models using Normalizing Flow

# 2. Models using Normalizing Flow

# 2. Models using Normalizing Flow

## 2-1. Real NVP

**Real-valued Non-Volume Preserving ( Dinh et al., 2017 )**

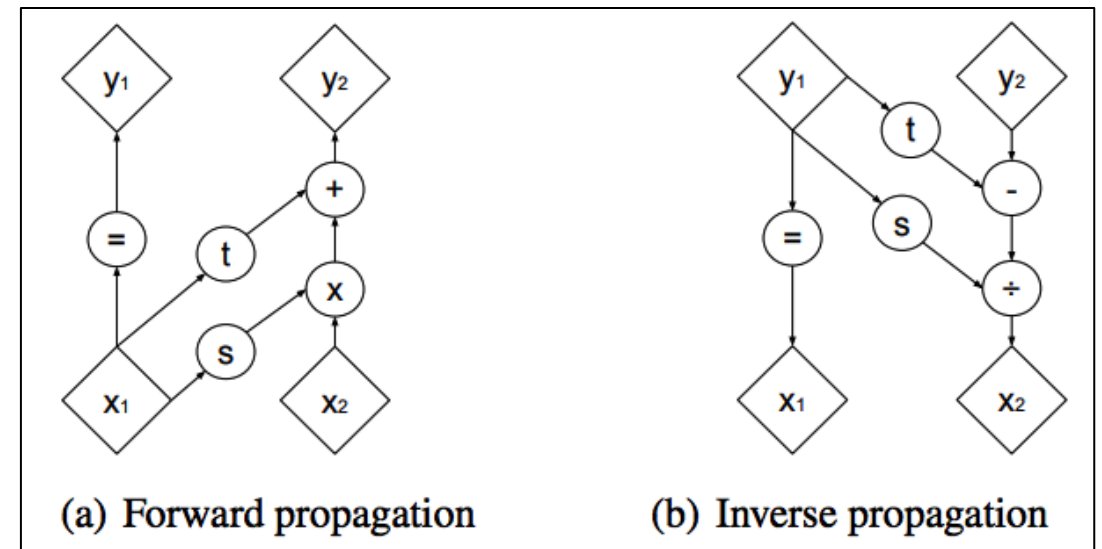Key point : stack multiple invertible bijective transformation, "Affine Coupling Layer"

Affine Coupling Layer

$$\mathbf{y}_{1:d} = \mathbf{x}_{1:d}$$
$$\mathbf{y}_{d+1:D} = \mathbf{x}_{d+1:D} \odot \exp(s(\mathbf{x}_{1:d})) + t(\mathbf{x}_{1:d})$$

Split the dimensions into 2 parts

- one part : Transformation (O)

- second part : Transformation (X)



(a) Forward propagation    (b) Inverse propagation

# 2. Models using Normalizing Flow

## 2-1. Real NVP

**Real-valued Non-Volume Preserving ( Dinh et al., 2017 )**

1. Easily Invertible

$$
\begin{cases} \mathbf{y}_{1:d} & = \mathbf{x}_{1:d} \\ \mathbf{y}_{d+1:D} & = \mathbf{x}_{d+1:D} \odot \exp(s(\mathbf{x}_{1:d})) + t(\mathbf{x}_{1:d}) \end{cases} \Leftrightarrow \begin{cases} \mathbf{x}_{1:d} & = \mathbf{y}_{1:d} \\ \mathbf{x}_{d+1:D} & = (\mathbf{y}_{d+1:D} - t(\mathbf{y}_{1:d})) \odot \exp(-s(\mathbf{y}_{1:d})) \end{cases}
$$

2. Determinant of Jacobian to be easily computed

$$
\frac{\partial y}{\partial x^T} = \begin{bmatrix} \mathbb{I}_d & 0 \\ \frac{\partial y_{d+1:D}}{\partial x_{1:d}^T} & \mathrm{diag}(\exp[s(x_{1:d})]) \end{bmatrix} \longrightarrow \exp\left[\sum_j s(x_{1:d})_j\right]
$$

## 2-2. NICE

**Non-linear Independent Component Estimation ( Dinh et al., 2015 )**

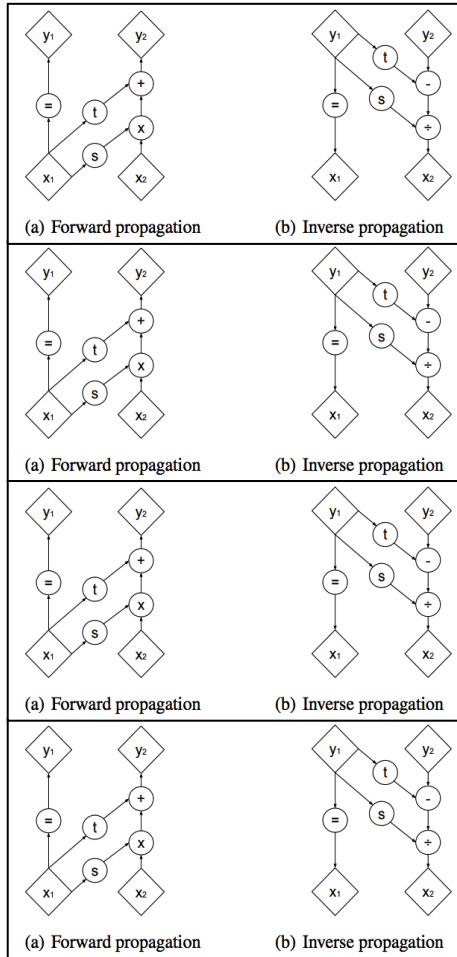Just make the Affine Coupling layer (of RealNVP) -> Additive Coupling layer

$$\begin{cases} \mathbf{y}_{1:d} & = \mathbf{x}_{1:d} \\ \mathbf{y}_{d+1:D} & = \mathbf{x}_{d+1:D} \odot \exp(s(\mathbf{x}_{1:d})) + t(\mathbf{x}_{1:d}) \end{cases} \Leftrightarrow \begin{cases} \mathbf{x}_{1:d} & = \mathbf{y}_{1:d} \\ \mathbf{x}_{d+1:D} & = (\mathbf{y}_{d+1:D} - t(\mathbf{y}_{1:d})) \odot \exp(-s(\mathbf{y}_{1:d})) \end{cases}$$

Affine Coupling Layer of RealNVP

$$\begin{cases} \mathbf{y}_{1:d} & = \mathbf{x}_{1:d} \\ \mathbf{y}_{d+1:D} & = \mathbf{x}_{d+1:D} + m(\mathbf{x}_{1:d}) \end{cases} \Leftrightarrow \begin{cases} \mathbf{x}_{1:d} & = \mathbf{y}_{1:d} \\ \mathbf{x}_{d+1:D} & = \mathbf{y}_{d+1:D} - m(\mathbf{y}_{1:d}) \end{cases}$$

Additive Coupling Layer of NICE

(a) Forward propagation  (b) Inverse propagation

(a) Forward propagation  (b) Inverse propagation

(a) Forward propagation  (b) Inverse propagation

(a) Forward propagation  (b) Inverse propagation

Think in advance. What would be the problem,

if we just stack the layers of RealNVP, NICE?

Think in advance. What would be the problem,

if we just stack the layers of RealNVP, NICE?
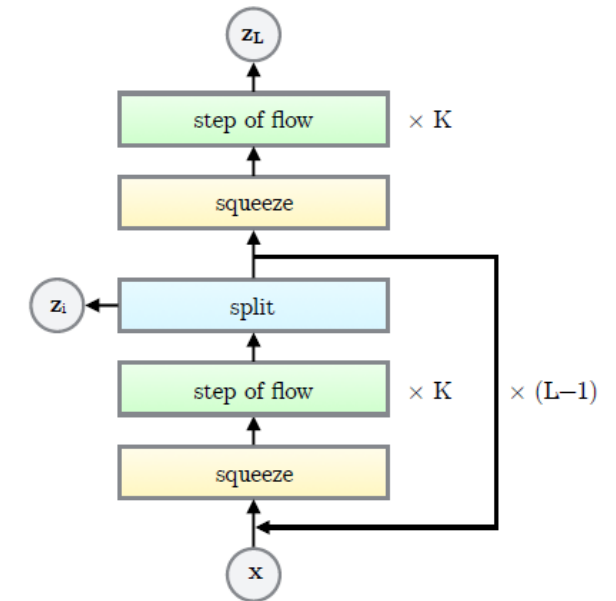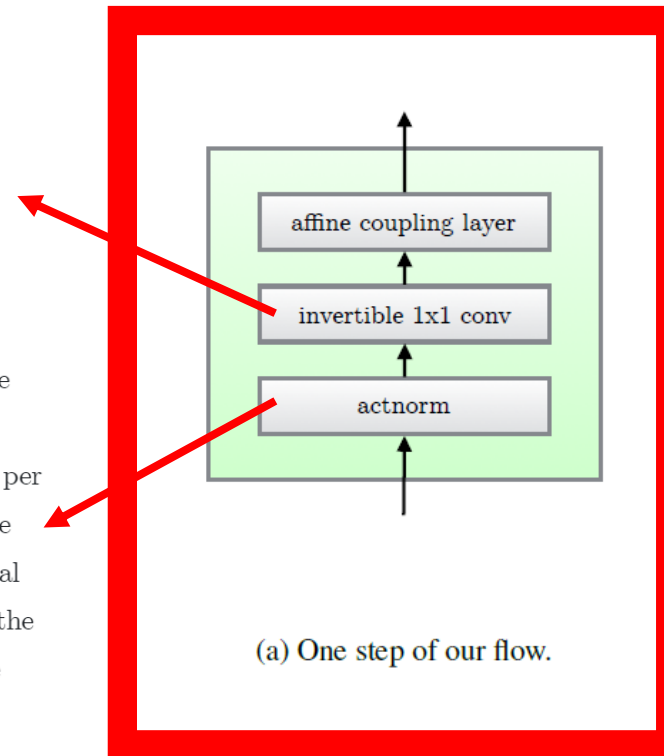
**Need permutation between the dimensions!**

## 2-3. Glow

**Generative Flow ( Kingma and Dhariwal, 2018 )**

Generalization of Permutation!

Activation Normalization is a type of normalization used for flow-based generative models; specifically it was introduced in the GLOW architecture. An ActNorm layer performs an affine transformation of the activations using a scale and bias parameter per channel, similar to batch normalization. These parameters are initialized such that the post-actnorm activations per-channel have zero mean and unit variance given an initial minibatch of data. This is a form of data dependent initilization. After initialization, the scale and bias are treated as regular trainable parameters that are independent of the data.

affine coupling layer

invertible 1x1 conv

actnorm

(a) One step of our flow.

$z_L$

step of flow × K

squeeze

$z_i$ ← split

step of flow × K × (L−1)

squeeze

x

(b) Multi-scale architecture (Dinh et al., 2016).

## 2-3. Glow

**Generative Flow ( Kingma and Dhariwal, 2018 )**

1x1 Convolution

$h \times w \times c$ tensor h with $c \times c$ weight matrix $\mathbf{W}$

$$\log \left| \det \left( \frac{d\, \mathbf{conv2D}(\mathbf{h}; \mathbf{W})}{d\, \mathbf{h}} \right) \right| = h \cdot w \cdot \log |\det(\mathbf{W})|$$

➡ cost of computing or differentiating $\det(\mathbf{W})$ is $\mathcal{O}\left(c^3\right)$

**How to solve this impracticability??**

## 2-3. Glow

**Generative Flow ( Kingma and Dhariwal, 2018 )**

1x1 Convolution

Use **LU Decomposition :** $\mathcal{O}\left(c^3\right) \rightarrow \mathcal{O}(c)$

$$\mathbf{W} = \mathbf{PL}(\mathbf{U} + \mathrm{diag}(\mathbf{s}))$$

- P : permutation matrix
- L is a lower triangular matrix with ones on the diagonal
- U is an upper triangular matrix with zeros on the diagonal

$$\log|\det(\mathbf{W})| = \mathrm{sum}(\log|\mathbf{s}|)$$

# 2. Models using Normalizing Flow

## 2-3. Glow

**Generative Flow ( Kingma and Dhariwal, 2018 )**

| Description | Function | Reverse Function | Log-determinant |
|---|---|---|---|
| Actnorm. See Section 3.1. | $\forall i, j : \mathbf{y}_{i,j} = \mathbf{s} \odot \mathbf{x}_{i,j} + \mathbf{b}$ | $\forall i, j : \mathbf{x}_{i,j} = (\mathbf{y}_{i,j} - \mathbf{b})/\mathbf{s}$ | $h \cdot w \cdot \text{sum}(\log|\mathbf{s}|)$ |
| Invertible $1 \times 1$ convolution. $\mathbf{W} : [c \times c]$. See Section 3.2. | $\forall i, j : \mathbf{y}_{i,j} = \mathbf{W}\mathbf{x}_{i,j}$ | $\forall i, j : \mathbf{x}_{i,j} = \mathbf{W}^{-1}\mathbf{y}_{i,j}$ | $h \cdot w \cdot \log|\det(\mathbf{W})|$ or $h \cdot w \cdot \text{sum}(\log|\mathbf{s}|)$ (see eq. (10)) |
| Affine coupling layer. See Section 3.3 and (Dinh et al., 2014) | $\mathbf{x}_a, \mathbf{x}_b = \texttt{split}(\mathbf{x})$ $(\log\mathbf{s}, \mathbf{t}) = \texttt{NN}(\mathbf{x}_b)$ $\mathbf{s} = \exp(\log\mathbf{s})$ $\mathbf{y}_a = \mathbf{s} \odot \mathbf{x}_a + \mathbf{t}$ $\mathbf{y}_b = \mathbf{x}_b$ $\mathbf{y} = \texttt{concat}(\mathbf{y}_a, \mathbf{y}_b)$ | $\mathbf{y}_a, \mathbf{y}_b = \texttt{split}(\mathbf{y})$ $(\log\mathbf{s}, \mathbf{t}) = \texttt{NN}(\mathbf{y}_b)$ $\mathbf{s} = \exp(\log\mathbf{s})$ $\mathbf{x}_a = (\mathbf{y}_a - \mathbf{t})/\mathbf{s}$ $\mathbf{x}_b = \mathbf{y}_b$ $\mathbf{x} = \texttt{concat}(\mathbf{x}_a, \mathbf{x}_b)$ | $\text{sum}(\log(|\mathbf{s}|))$ |

Fig. 4. Three substeps in one step of flow in Glow. (Image source: Kingma and Dhariwal, 2018)

# 3. Models using Autoregressive Flow

# 3. Models using Autoregressive Flow

Autoregressive = only depends on the "previous" parts

$$p(\mathbf{x}) = \prod_{i=1}^{D} p\left(x_i \mid x_1, \ldots, x_{i-1}\right) = \prod_{i=1}^{D} p\left(x_i \mid x_{1:i-1}\right)$$

# 3. Models using Autoregressive Flow

Autoregressive = only depends on the "previous" parts

$$p(\mathbf{x}) = \prod_{i=1}^{D} p(x_i \mid x_1, \ldots, x_{i-1}) = \prod_{i=1}^{D} p(x_i \mid x_{1:i-1})$$

This structure enables "triangular Jacobian"

3 popular model using autoregressive structure

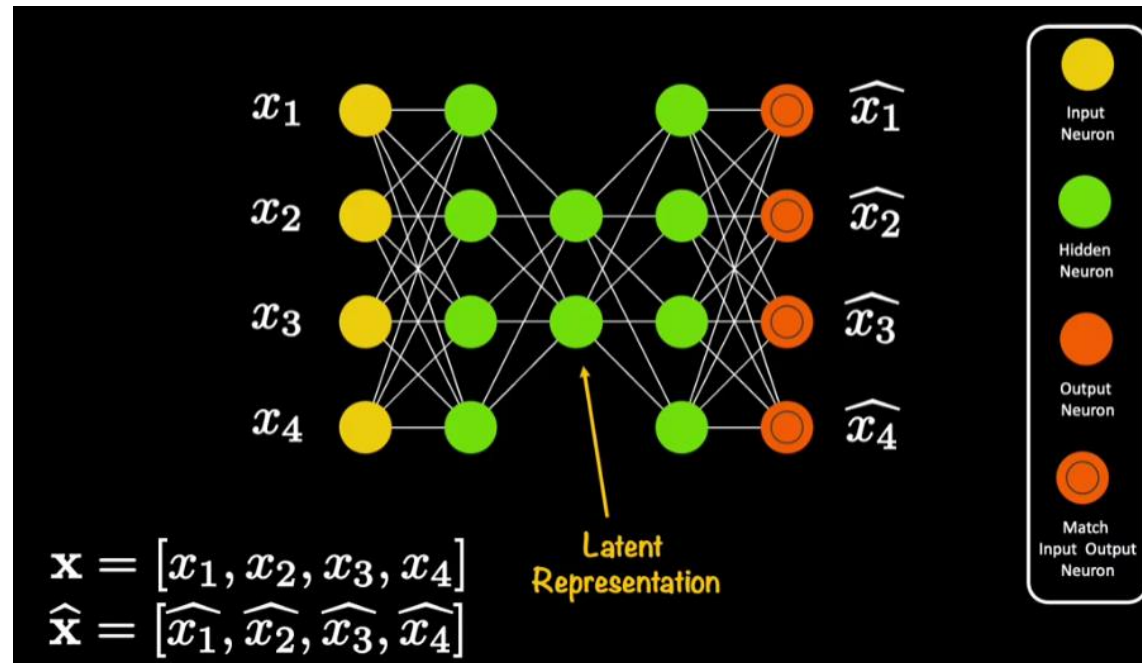    **1) MADE ( Masked Autoencoder Density Estimation )**

    ( Autoregressive Flow )

    **2) MAF ( Masked Autoregressive Flow )**

    **3) IAF ( Inverse Autoregressive Flow )**

## 3-1. MADE

**Masked Autoencoder Density Estimation ( Germain et al., 2015 )**



https://www.youtube.com/watch?v=7q4ueFiJjAY
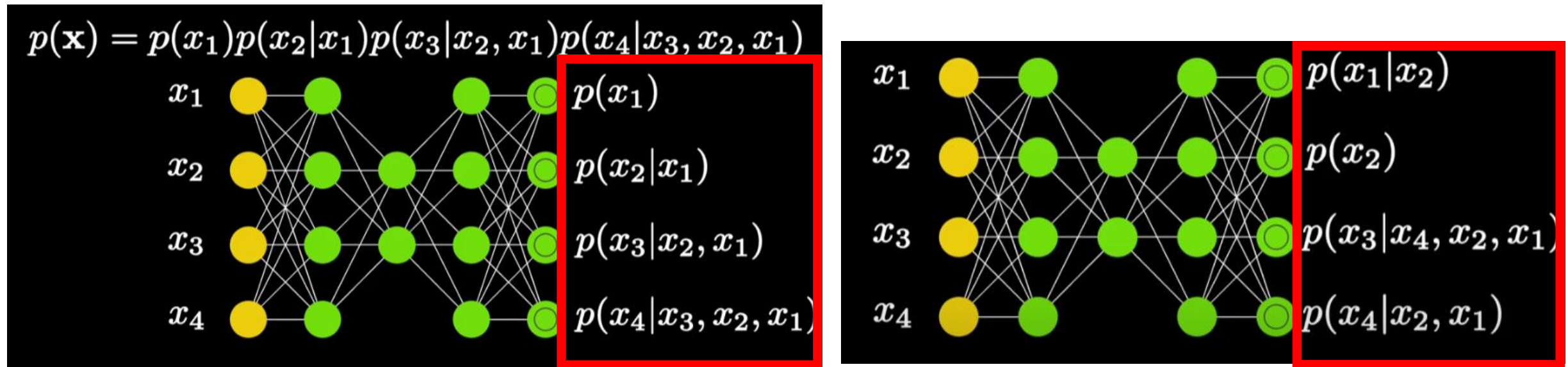
## 3-1. MADE

**Masked Autoencoder Density Estimation ( Germain et al., 2015 )**



$$p(\mathbf{x}) = p(x_1)p(x_2)p(x_3)p(x_4)$$

No conditional dependence between the input variables!

## 3-1. MADE

**Masked Autoencoder Density Estimation ( Germain et al., 2015 )**



conditional dependence (O)  **( Output Ordering is arbitrary )**

# 3. Models using Autoregressive Flow

## 3-1. MADE

**Masked Autoencoder Density Estimation ( Germain et al., 2015 )**

Goal : Constrain the AE to learn joint pdf with conditional dependence!
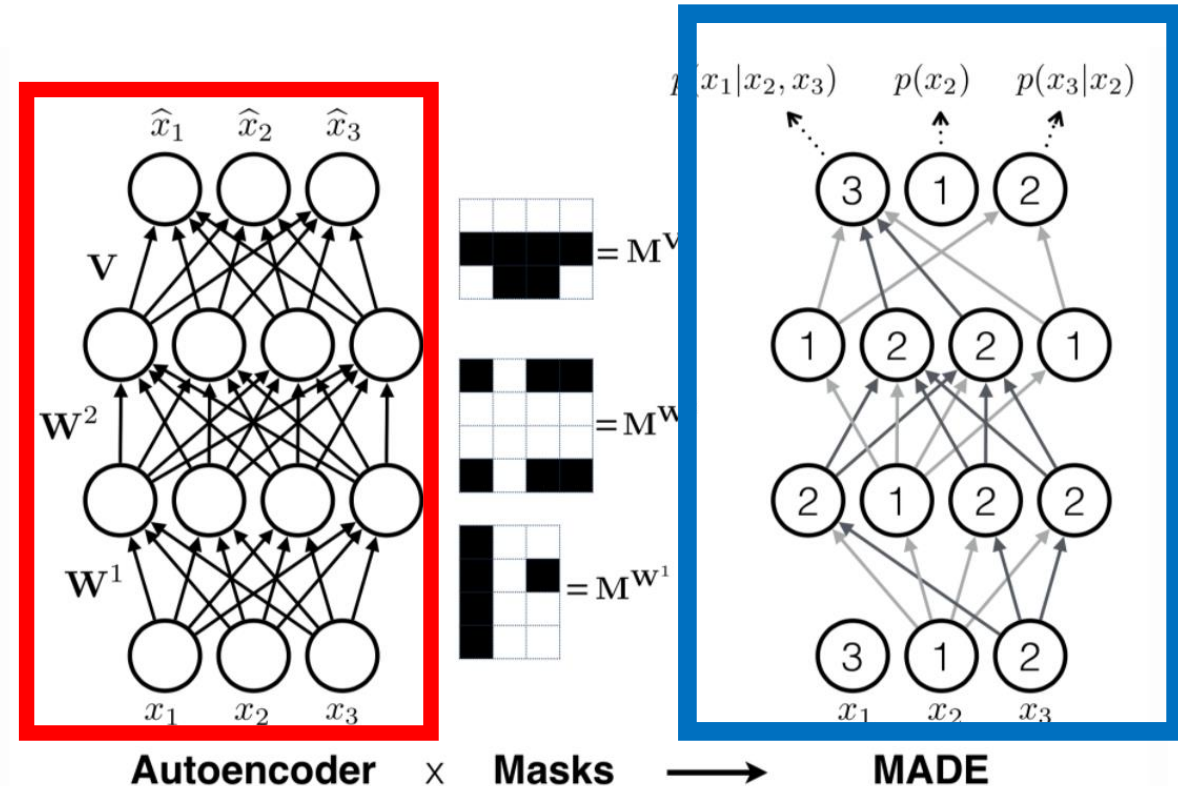
(+agnostic dimension ordering)

## 3-1. MADE

**Masked Autoencoder Density Estimation ( Germain et al., 2015 )**

Way to disconnect the weights : use binary Masks!

$$\mathbf{h}^0 = \mathbf{x}$$
$$\mathbf{h}^l = \text{activation}^l(\mathbf{W}^l \mathbf{h}^{l-1} + \mathbf{b}^l)$$
$$\hat{\mathbf{x}} = \sigma(\mathbf{V}\mathbf{h}^L + \mathbf{c})$$

$$\mathbf{h}^l = \text{activation}^l((\mathbf{W}^l \odot \mathbf{M}^{\mathbf{W}^l})\mathbf{h}^{l-1} + \mathbf{b}^l)$$
$$\hat{\mathbf{x}} = \sigma((\mathbf{V} \odot \mathbf{M}^{\mathbf{V}})\mathbf{h}^L + \mathbf{c})$$
$$M^{\mathbf{W}^l}_{k',k} = \mathbf{1}_{m^l_{k'} \geq m^{l-1}_k} = \begin{cases} 1, & \text{if } m^l_{k'} \geq m^{l-1}_k \\ 0, & \text{otherwise} \end{cases}$$
$$M^{\mathbf{V}}_{d,k} = \mathbf{1}_{d \geq m^L_k} = \begin{cases} 1, & \text{if } d > m^L_k \\ 0, & \text{otherwise} \end{cases}$$



Autoencoder x Masks ⟶ MADE

## 3-1. MADE

**Masked Autoencoder Density Estimation ( Germain et al., 2015 )**

Then, how can we impose "Autoregressive" property to mask?

- 1) assign each unit (in hidden layer) an integer $m$ ( between $1$ and $D-1$ )

  ( $m(k)$ = maximum number of input units to which it can be connected )

  ( $m(k) \neq 1, m(k) \neq D$ )

Notation & Meaning

- $\mathbf{M}^V$ and $\mathbf{M}^W$ : network's connectivity
- matrix product $\mathbf{M}^{V,W} = \mathbf{M}^V \mathbf{M}^W$ : connectivity between the input and the output layer
- $M_{d',d}^{V,W}$ : number of network paths between output unit $\hat{x}_{d'}$ and input unit $x_d$.

- 2) [MASK] matrix masking the connections between "input & hidden units"

  constraints on the maximum number of inputs to each hidden unit are encoded in it!

$$M_{k,d}^{W} = 1_{m(k) \geq d} = \begin{cases} 1 & \text{if } m(k) \geq d \\ 0 & \text{otherwise} \end{cases}$$

- 3) [MASK] matrix masking the connections between "hidden & output" units

$$M_{d,k}^{V} = 1_{d > m(k)} = \begin{cases} 1 & \text{if } d > m(k) \\ 0 & \text{otherwise} \end{cases}$$

**( single hidden layer NN )**

# 3. Models using Autoregressive Flow

## 3-1. MADE

**Masked Autoencoder Density Estimation ( Germain et al., 2015 )**

Then, how can we impose "Autoregressive" property to mask?

Mask

- $M_{k',k}^{\mathbf{W}^l} = 1_{m^l(k') \geq m^{l-1}(k)} = \begin{cases} 1 & \text{if } m^l(k') \geq m^{l-1}(k) \\ 0 & \text{otherwise} \end{cases}$

- $M_{d,k}^{\mathbf{V}} = 1_{d > m^L(k)} = \begin{cases} 1 & \text{if } d > m^L(k) \\ 0 & \text{otherwise} \end{cases}$
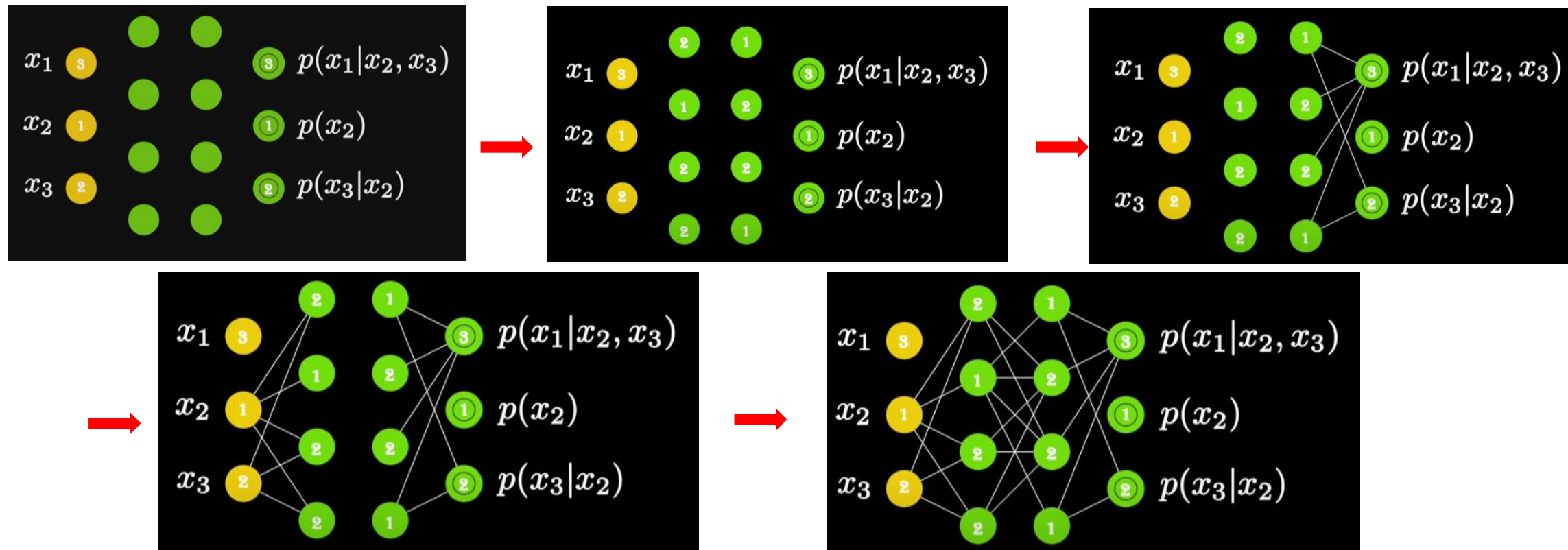
Notation

- $\mathbf{W}^1$ : first hidden layer matrix
- $\mathbf{W}^2$ : second hidden layer matrix
- $K^l$ : number of hidden units in layer $l$
- $m^l(k)$ : maximum number of connected inputs of the $k^{\text{th}}$ unit in the $l^{\text{th}}$ layer        **( Deep NN )**

# 3. Models using Autoregressive Flow

## 3-1. MADE

**Masked Autoencoder Density Estimation ( Germain et al., 2015 )**

## 3-1. MADE

**Masked Autoencoder Density Estimation ( Germain et al., 2015 )**

## 3-1. MADE

**Masked Auto**

Autoregressive Flow

$$x_1 \sim p_\theta\left(x_1\right)$$

$$x_2 \sim p_\theta\left(x_2 \mid x_1\right)$$

$$x_3 \sim p_\theta\left(x_3 \mid x_1, x_2\right)$$

Autoregressive Flow

$$x_1 \sim p_\theta(x_1)$$

$$x_2 \sim p_\theta(x_2 \mid x_1)$$

$$x_3 \sim p_\theta(x_3 \mid x_1, x_2)$$

$$x_1 = f_\theta^{-1}(z_1)$$

$$x_2 = f_\theta^{-1}(z_2; x_1)$$

$$x_3 = f_\theta^{-1}(z_3; x_1, x_2)$$

$$z_i \sim N(0, 1)$$

**1) Sampling**

# 3. Models using Autoregressive Flow

Seunghan Lee, Yonsei University

**Autoregressive Flow**

$$x_1 \sim p_\theta(x_1)$$

$$x_2 \sim p_\theta(x_2 \mid x_1)$$

$$x_3 \sim p_\theta(x_3 \mid x_1, x_2)$$

$$x_1 = f_\theta^{-1}(z_1)$$

$$x_2 = f_\theta^{-1}(z_2; x_1)$$

$$x_3 = f_\theta^{-1}(z_3; x_1, x_2)$$

$$z_i \sim N(0, 1)$$

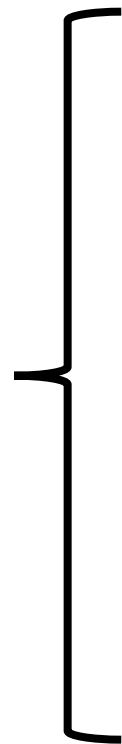**1) Sampling**   **Can be done in "D" pass (Slow)**

# 3. Models using Autoregressive Flow

Autoregressive Flow

$$x_1 \sim p_\theta(x_1)$$

$$x_2 \sim p_\theta(x_2 \mid x_1)$$

$$x_3 \sim p_\theta(x_3 \mid x_1, x_2)$$

$$z_i \sim N(0,1)$$

$$x_1 = f_\theta^{-1}(z_1)$$

$$x_2 = f_\theta^{-1}(z_2; x_1)$$

$$x_3 = f_\theta^{-1}(z_3; x_1, x_2)$$

**1) Sampling**     **Can be done in "D" pass (Slow)**

$$z_1 = f_\theta(x_1)$$

$$z_2 = f_\theta(x_2; x_1)$$

$$z_3 = f_\theta(x_3; x_1, x_2)$$

**2) Density Evaluation**

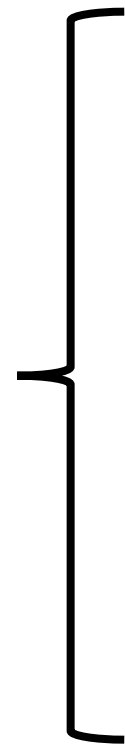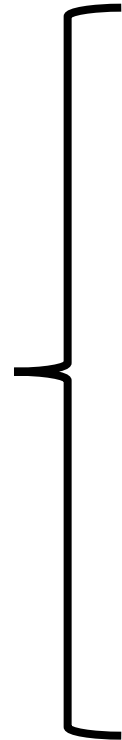Seunghan Lee, Yonsei University

# 3. Models using Autoregressive Flow

Autoregressive Flow

$$x_1 \sim p_\theta(x_1)$$

$$x_2 \sim p_\theta(x_2 \mid x_1)$$

$$x_3 \sim p_\theta(x_3 \mid x_1, x_2)$$

$$x_1 = f_\theta^{-1}(z_1)$$

$$x_2 = f_\theta^{-1}(z_2; x_1)$$

$$x_3 = f_\theta^{-1}(z_3; x_1, x_2)$$

$$z_i \sim N(0,1)$$

**1) Sampling**  **Can be done in "D" pass (Slow)**

$$z_1 = f_\theta(x_1)$$

$$z_2 = f_\theta(x_2; x_1)$$

$$z_3 = f_\theta(x_3; x_1, x_2)$$

**2) Density Evaluation**  **Can be done in 1 pass (Fast)**

# 3. Models using Autoregressive Flow

## 3-2. MAF

**Masked Autoregressive Flow ( Papamakarios et al, 2017 )**

**MAF = Flow version of MADE**

It also follows an autoregressive property

$$p(\mathbf{x}) = \prod_{i=1}^{D} p(x_i | \mathbf{x}_{1:i-1})$$

# 3. Models using Autoregressive Flow

## 3-2. MAF

**Masked Autoregressive Flow ( Papamakarios et al, 2017 )**

**MAF = Flow version of MADE**

It also follows an autoregressive property $\quad p(\mathbf{x}) = \prod_{i=1}^{D} p(x_i | \mathbf{x}_{1:i-1})$

Autoregressive model whose conditionals are parameterized as **single Gaussians**

$$p\left(x_i \mid \mathbf{x}_{1:i-1}\right) = \mathcal{N}\left(x_i \mid \mu_i, (\exp \alpha_i)^2\right)$$

$$\text{where} \quad \mu_i = f_{\mu_i}\left(\mathbf{x}_{1:i-1}\right) \text{ and } \alpha_i = f_{\alpha_i}\left(\mathbf{x}_{1:i-1}\right)$$

## 3-2. MAF

**Masked Autoregressive Flow ( Papamakarios et al, 2017 )**

**MAF = Flow version of MADE**

It also follows an autoregressive property $\quad p(\mathbf{x}) = \prod_{i=1}^{D} p(x_i | \mathbf{x}_{1:i-1})$

Autoregressive model whose conditionals are parameterized as **single Gaussians**

$$p\left(x_i \mid \mathbf{x}_{1:i-1}\right) = \mathcal{N}\left(x_i \mid \mu_i, (\exp \alpha_i)^2\right)$$

$$\text{where} \quad \mu_i = f_{\mu_i}\left(\mathbf{x}_{1:i-1}\right) \text{ and } \alpha_i = f_{\alpha_i}\left(\mathbf{x}_{1:i-1}\right)$$

use "MADE"

## 3-2. MAF

**Masked Autoregressive Flow ( Papamakarios et al, 2017 )**

**MAF = Flow version of MADE**

It also follows an autoregressive property $\quad p(\mathbf{x}) = \prod_{i=1}^{D} p(x_i | \mathbf{x}_{1:i-1})$

Autoregressive model whose conditionals are parameterized as **single Gaussians**

$$p\left(x_i \mid \mathbf{x}_{1:i-1}\right) = \mathcal{N}\left(x_i \mid \mu_i, (\exp \alpha_i)^2\right)$$

where $\quad \mu_i = f_{\mu_i}\left(\mathbf{x}_{1:i-1}\right)$ and $\alpha_i = f_{\alpha_i}\left(\mathbf{x}_{1:i-1}\right)$

use "MADE"

$\Longrightarrow$

$$x_i = u_i \exp \alpha_i + \mu_i$$
- $\mu_i = f_{\mu_i}\left(\mathbf{x}_{1:i-1}\right)$
- $\alpha_i = f_{\alpha_i}\left(\mathbf{x}_{1:i-1}\right)$
- $u_i \sim \mathcal{N}(0,1)$

# 3. Models using Autoregressive Flow

## 3-2. MAF

**Masked Autoregressive Flow ( Papamakarios et al, 2017 )**

**MAF = Flow version of MADE**

It also follows an autoregressive property $\quad p(\mathbf{x}) = \prod_{i=1}^{D} p(x_i | \mathbf{x}_{1:i-1})$

(SAMPLING) We can generate data, using recursion :

$x_i = u_i \exp \alpha_i + \mu_i$

- $\mu_i = f_{\mu_i}(\mathbf{x}_{1:i-1})$
- $\alpha_i = f_{\alpha_i}(\mathbf{x}_{1:i-1})$
- $u_i \sim \mathcal{N}(0,1)$

$\blacktriangleright$

1. $x_1 = f_{\mu_1} + \exp(f_{\sigma_1})z_1$ for $z_1 \sim N(0,1)$
2. $x_2 = f_{\mu_2}(x_1) + \exp(f_{\sigma_2}(x_1))z_2$ for $z_2 \sim N(0,1)$
3. $x_3 = f_{\mu_3}(x_1, x_2) + \exp(f_{\sigma_3}(x_1, x_2))z_3$ for $z_3 \sim N(0,1)$

# 3. Models using Autoregressive Flow

Inverse Autoregressive Flow

$$z_1 \sim p_\theta\left(z_1\right)$$

$$z_2 \sim p_\theta\left(z_2 \mid z_1\right)$$

$$z_3 \sim p_\theta\left(z_3 \mid z_1, z_2\right)$$

| | |
|---|---|
| | $x_1 \sim p_\theta\left(x_1\right)$ |
| | $x_2 \sim p_\theta\left(x_2 \mid x_1\right)$ |
| ( Autoregressive Flow ) | $x_3 \sim p_\theta\left(x_3 \mid x_1, x_2\right)$ |

Inverse Autoregressive Flow

$$z_i \sim N(0, 1)$$

$$x_1 = f_\theta(z_1)$$
$$x_2 = f_\theta(z_2; z_1)$$
$$x_3 = f_\theta(z_3; z_1, z_2)$$

**1) Sampling**

$$z_1 \sim p_\theta(z_1)$$
$$z_2 \sim p_\theta(z_2 \mid z_1)$$
$$z_3 \sim p_\theta(z_3 \mid z_1, z_2)$$

$$x_1 \sim p_\theta(x_1)$$
$$x_2 \sim p_\theta(x_2 \mid x_1)$$
( Autoregressive Flow ) $\quad x_3 \sim p_\theta(x_3 \mid x_1, x_2)$

Inverse Autoregressive Flow

$$z_i \sim N(0,1)$$

$$x_1 = f_\theta(z_1)$$
$$x_2 = f_\theta(z_2; z_1)$$
$$x_3 = f_\theta(z_3; z_1, z_2)$$

**1) Sampling**  **Can be done in 1 pass (Fast)**

$$z_1 \sim p_\theta(z_1)$$
$$z_2 \sim p_\theta(z_2 \mid z_1)$$
$$z_3 \sim p_\theta(z_3 \mid z_1, z_2)$$

$$x_1 \sim p_\theta(x_1)$$
$$x_2 \sim p_\theta(x_2 \mid x_1)$$
( Autoregressive Flow )   $$x_3 \sim p_\theta(x_3 \mid x_1, x_2)$$

# 3. Models using Autoregressive Flow

**Inverse Autoregressive Flow**

$$z_i \sim N(0, 1)$$

$$x_1 = f_\theta(z_1)$$
$$x_2 = f_\theta(z_2; z_1)$$
$$x_3 = f_\theta(z_3; z_1, z_2)$$

**1) Sampling**    **Can be done in 1 pass (Fast)**

$$z_1 \sim p_\theta(z_1)$$
$$z_2 \sim p_\theta(z_2 \mid z_1)$$
$$z_3 \sim p_\theta(z_3 \mid z_1, z_2)$$

$$z_1 = f_\theta^{-1}(x_1)$$
$$z_2 = f_\theta^{-1}(x_2; z_1)$$
$$z_3 = f_\theta^{-1}(x_3; z_1, z_2)$$

$$x_1 \sim p_\theta(x_1)$$
$$x_2 \sim p_\theta(x_2 \mid x_1)$$
( Autoregressive Flow )    $$x_3 \sim p_\theta(x_3 \mid x_1, x_2)$$

**2) Density Evaluation**

# 3. Models using Autoregressive Flow

**Inverse Autoregressive Flow**

$$z_1 \sim p_\theta\left(z_1\right)$$

$$z_2 \sim p_\theta\left(z_2 \mid z_1\right)$$

$$z_3 \sim p_\theta\left(z_3 \mid z_1, z_2\right)$$

$$x_1 \sim p_\theta\left(x_1\right)$$
$$x_2 \sim p_\theta\left(x_2 \mid x_1\right)$$
( Autoregressive Flow )  $\quad x_3 \sim p_\theta\left(x_3 \mid x_1, x_2\right)$

$$z_i \sim N(0, 1)$$

$$x_1 = f_\theta(z_1)$$
$$x_2 = f_\theta(z_2; z_1)$$
$$x_3 = f_\theta(z_3; z_1, z_2)$$

**1) Sampling** **Can be done in 1 pass (Fast)**

$$z_1 = f_\theta^{-1}(x_1)$$
$$z_2 = f_\theta^{-1}(x_2; z_1)$$
$$z_3 = f_\theta^{-1}(x_3; z_1, z_2)$$

**2) Density Evaluation** **Can be done in "D" pass (Slow)**

## 3-3. IAF

**Inverse Autoregressive Flow ( Kingma et al, 2016 )**

**IAF = Inverse version of MAF**

Instead of "X", **"U" passes through the flow!**

$$x_i = u_i \exp \alpha_i + \mu_i$$

- $u_i \sim N(0,1)$    $i = 1, \ldots, D$

< MAF >

$$u_i = \frac{(x_i - \mu_i)}{\exp \alpha_i}$$

- $u_1 \sim N(0,1)$

< IAF >

Seunghan Lee, Yonsei University

## 3-3. IAF

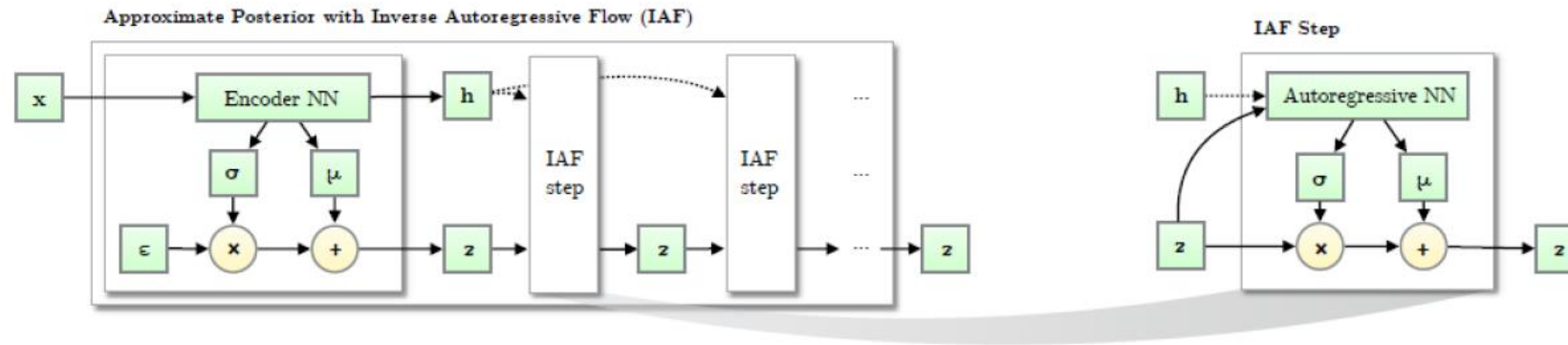**Inverse Autoregressive Flow ( Kingma et al, 2016 )**



Figure 2: Like other normalizing flows, drawing samples from an approximate posterior with Inverse Autoregressive Flow (IAF) consists of an initial sample **z** drawn from a simple distribution, such as a Gaussian with diagonal covariance, followed by a chain of nonlinear invertible transformations of **z**, each with a simple Jacobian determinants.

## 3-3. IAF

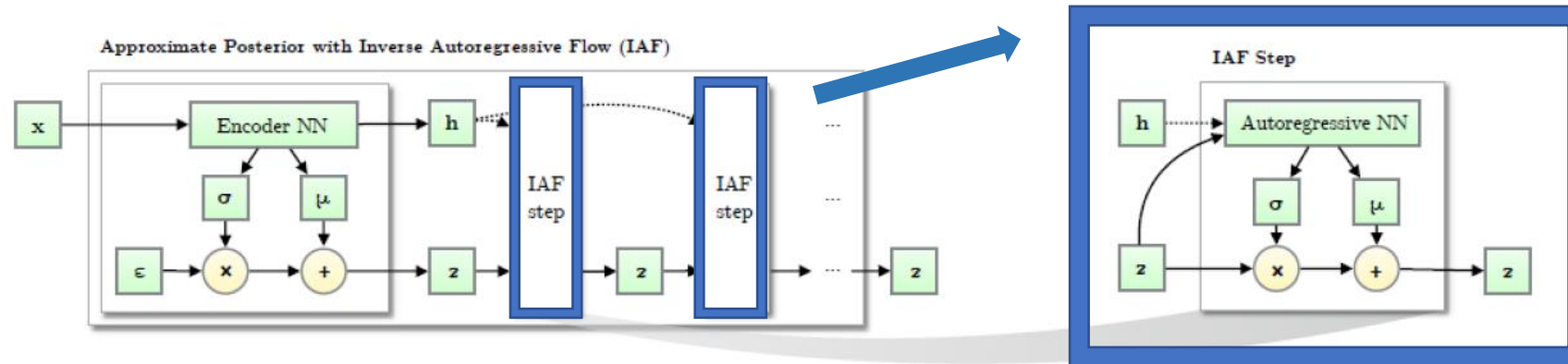**Inverse Autoregressive Flow ( Kingma et al, 2016 )**



Figure 2: Like other normalizing flows, drawing samples from an approximate posterior with Inverse Autoregressive Flow (IAF) consists of an initial sample **z** drawn from a simple distribution, such as a Gaussian with diagonal covariance, followed by a chain of nonlinear invertible transformations of **z**, each with a simple Jacobian determinants.

## 3-3. IAF

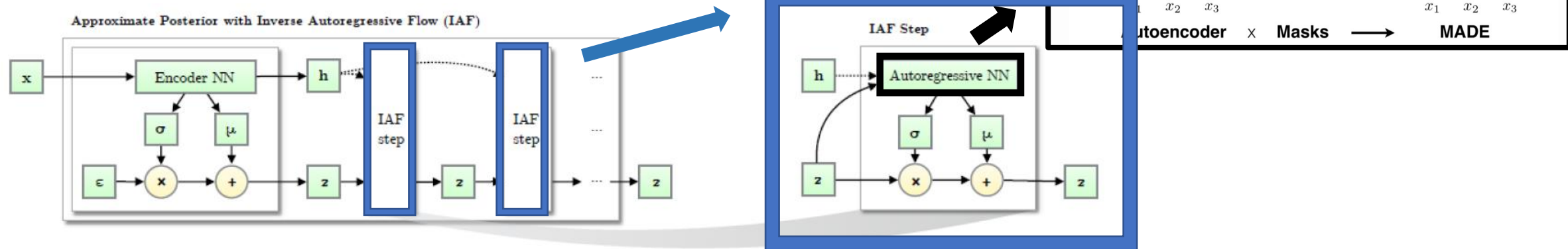**Inverse Autoregressive Flow ( Kingma et al, 2016 )**



Figure 2: Like other normalizing flows, drawing samples from an approximate posterior with Inverse Autoregressive Flow (IAF) consists of an initial sample **z** drawn from a simple distribution, such as a Gaussian with diagonal covariance, followed by a chain of nonlinear invertible transformations of **z**, each with a simple Jacobian determinants.

# 3. Models using Autoregressive Flow

## 3-3. IAF



**Algorithm 1:** Pseudo-code of an approximate posterior with Inverse Autoregressive Flow (IAF)

**Data:**

x: a datapoint, and optionally other conditioning information
$\theta$: neural network parameters
EncoderNN(x; $\theta$): encoder neural network, with additional output h
AutoregressiveNN[*](z, h; $\theta$): autoregressive neural networks, with additional input h
sum(.): sum over vector elements
sigmoid(.): element-wise sigmoid function

**Result:**

z: a random sample from $q(z|x)$, the approximate posterior distribution
l: the scalar value of $\log q(z|x)$, evaluated at sample 'z'

$[\mu, \sigma, h] \leftarrow \text{EncoderNN}(x; \theta)$
$\epsilon \sim \mathcal{N}(0, I)$
$z \leftarrow \sigma \odot \epsilon + \mu$
$l \leftarrow -\text{sum}(\log \sigma + \frac{1}{2}\epsilon^2 + \frac{1}{2}\log(2\pi))$
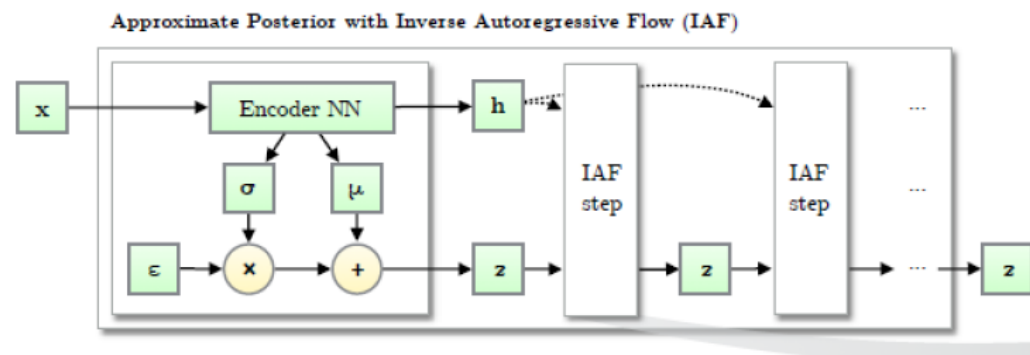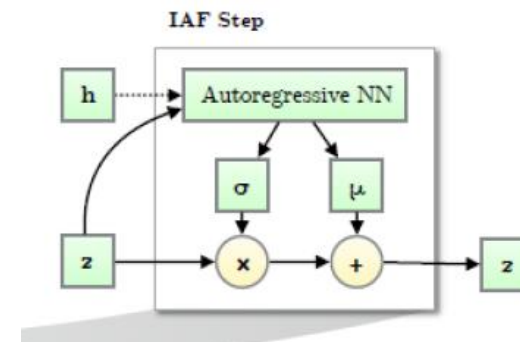**for** $t \leftarrow 1$ **to** $T$ **do**
$\quad [m, s] \leftarrow \text{AutoregressiveNN}[t](z, h; \theta)$
$\quad \sigma \leftarrow \text{sigmoid}(s)$
$\quad z \leftarrow \sigma \odot z + (1 - \sigma) \odot m$
$\quad l \leftarrow l - \text{sum}(\log \sigma)$
**end**

# 3. Models using Autoregressive Flow

## 3-3. IAF

**IAF vs MAF**

- [ MAF ]   FAST density evaluation, SLOW sampling
  - $\mu_i$ and $\alpha_i$ are directly computed from previous "data variables $x_{1:i-1}$"
  - capable of calculating the density $p(x)$ of any data point in one pass

    but sampling requires $D$ sequential passes

$$z_1 = f_\theta(x_1)$$
$$z_2 = f_\theta(x_2; x_1)$$
$$z_3 = f_\theta(x_3; x_1, x_2)$$

- [ IAF ]   FAST sampling, SLOW density evaluation
  - $\mu_i$ and $\alpha_i$ are directly computed from previous "random numbers $u_{1:i-1}$"
  - sampling requires only one pass

    but calculating the density $p(x)$ of any data point requires $D$ passes

$$x_1 = f_\theta(z_1)$$
$$x_2 = f_\theta(z_2; z_1)$$
$$x_3 = f_\theta(z_3; z_1, z_2)$$

# Summary

## Flow-based Model

- Change of variables

- Normalizing Flow

- Models using NF (Normalizing Flow)

    - Real NVP : affine coupling layer

    - NICE : additive coupling layer

    - Glow : 1x1 convolution & PLU decomposition

- Models using AF (Autoregressive Flow)

    - MADE

    - MAF : flow using MADE

    - IAF : inverse version of MAF

# Thank You

Seunghan Lee, Yonsei University