



Normalizing Flow

Il-Chul Moon

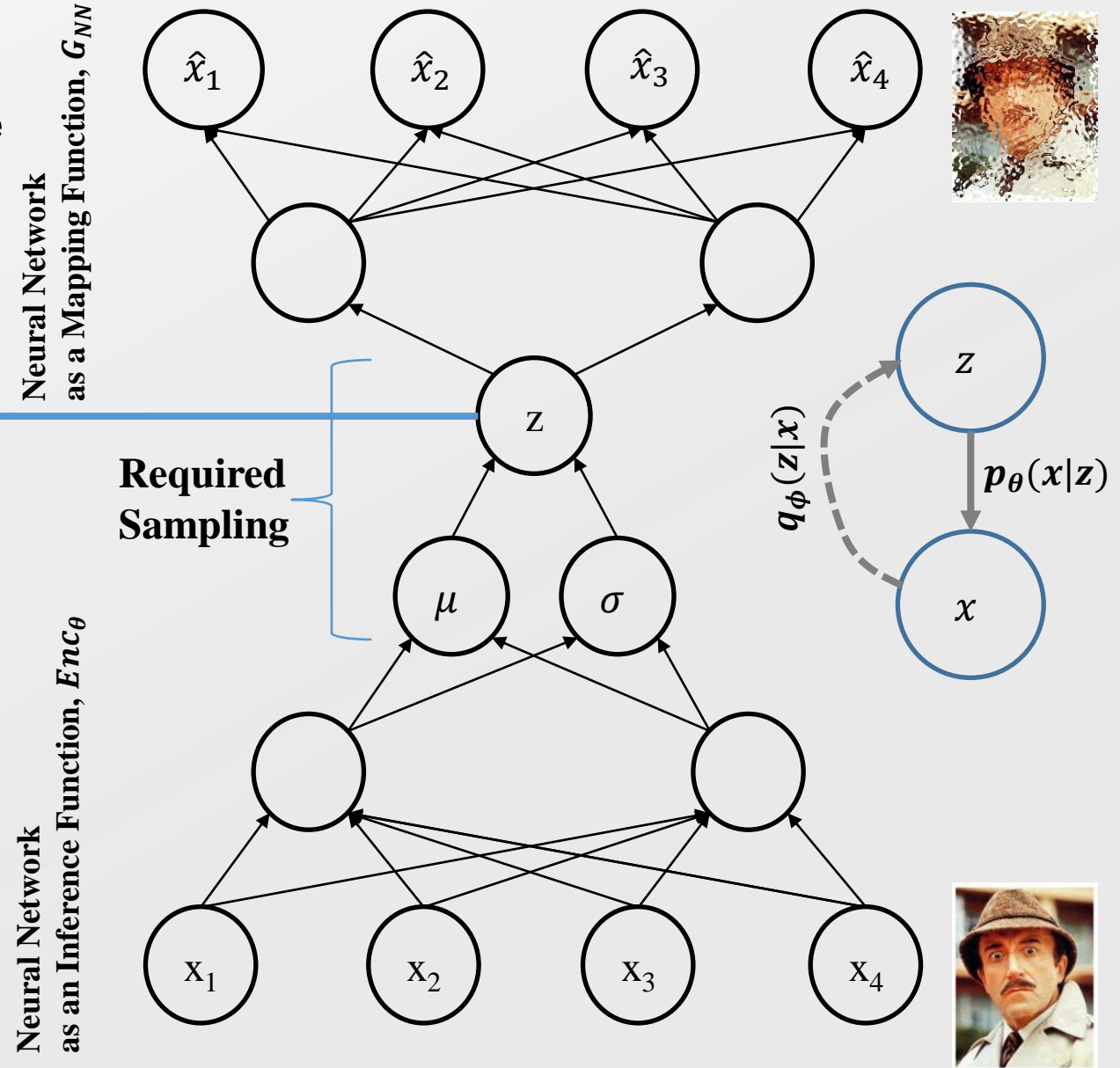
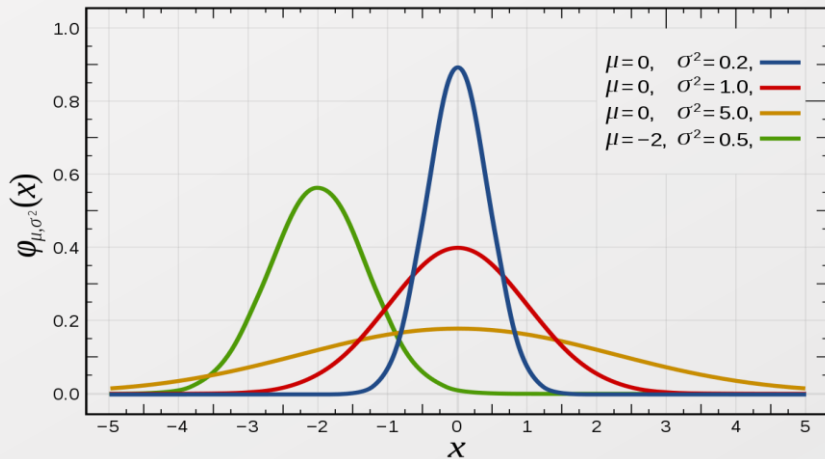
Department of Industrial and Systems Engineering

KAIST

icmoon@kaist.ac.kr

Detour: Explicit Deep Generative Model (VAE)

- Generative model + Neural network
 - Neural network is used for the inference of the distribution parameters
- $x_{data} \sim P_{data}, x_{sample} = \hat{x} \sim P_{gen}$
- $z \sim N(Enc_{\mu}(x_{data}), Enc_{\sigma}(x_{data}))$
- Minimizing
 - $Error(x, \hat{x}) = Error(x, G(z))$
 - with some regularizations
 - Prior on Z with $N(0, I)$

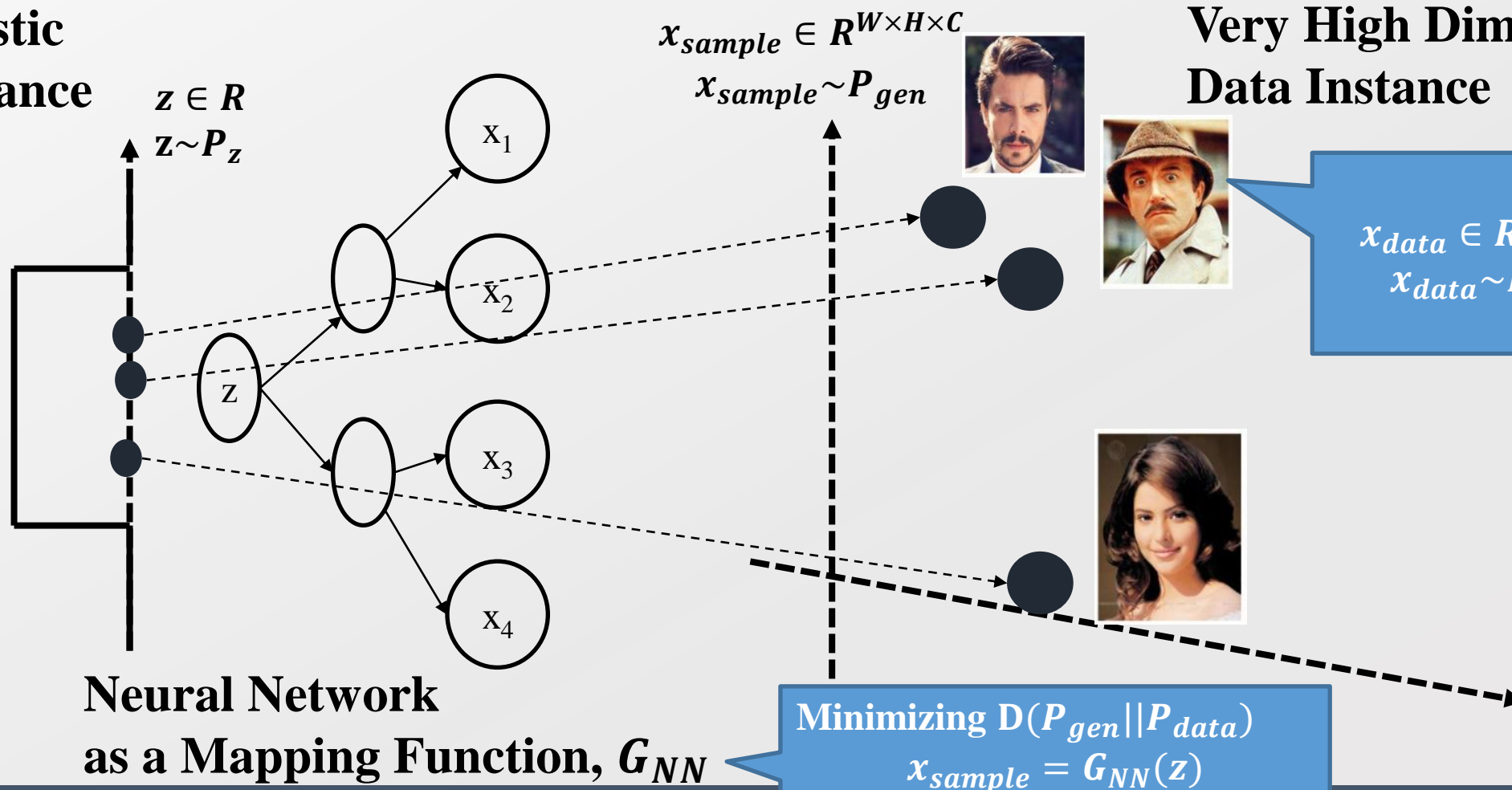


Detour: Implicit Deep Generative Model (GAN)

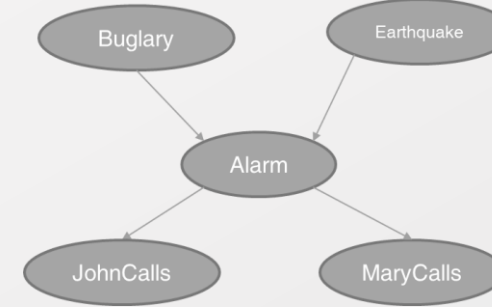
- Generative model + Neural network
 - Neural network is used for the specification of the density model

Function Approximation
on
Implicit Probability
Distribution of P_{gen}

**Low Stochastic
Sample Instance**

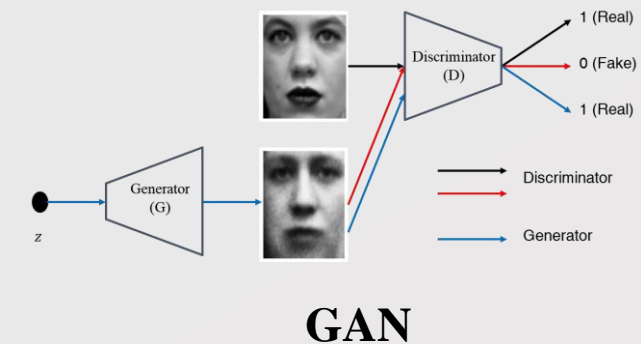
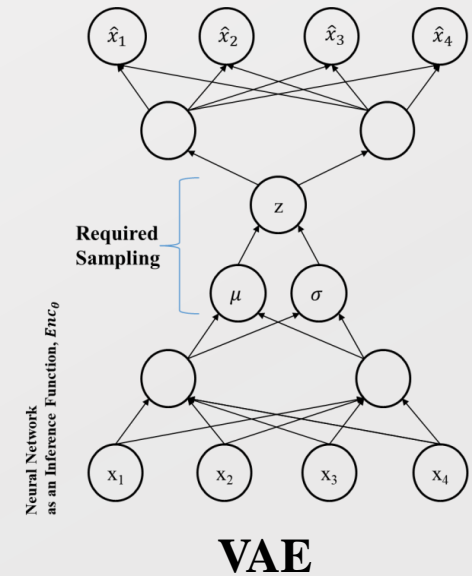
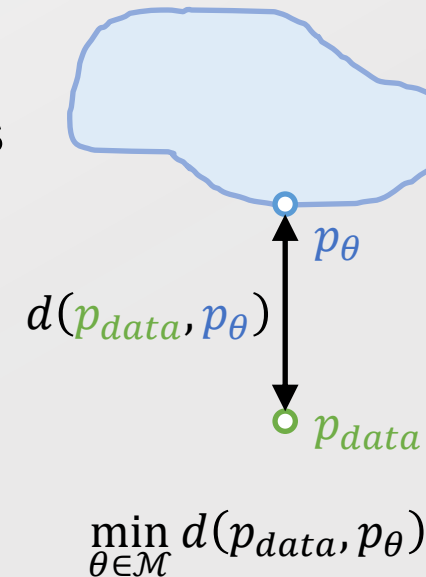


- Probabilistic Model, or Generative Models
 - Fundamental aspects
 - Generation (==Sampling, Simulation...)
 - Able to generate a sample from Bayesian Network, VAE, GAN
 - Inference
 - Able to calculate a likelihood from Bayesian Network and VAE
 - Inable to calculate a likelihood from GAN
 - a.k.a. Likelihood-Free Model, Likelihood-Free Inference
- Distinctions from the types of generative models
 - Given $x_j \sim p_{data}, j = 1, 2, \dots, |\mathcal{D}|$
 - 1) The flexibility of p_θ parameterized by $\theta \in \mathcal{M}$
 - 2) The property of $d(p_{data}, p_\theta)$
 - 3) The optimization of $\min_{\theta \in \mathcal{M}} d(p_{data}, p_\theta)$



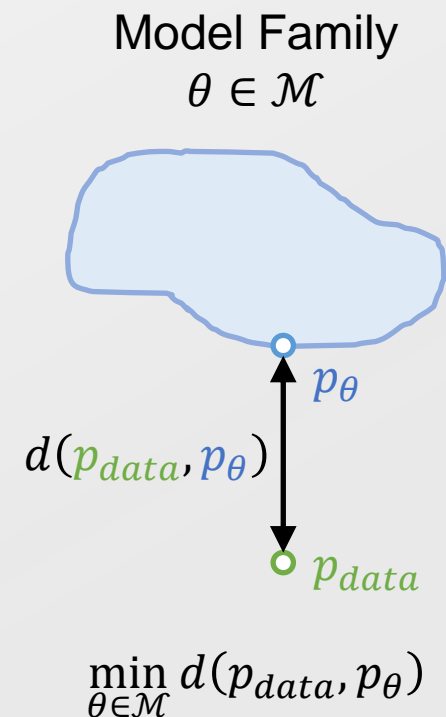
Bayesian Network

Model Family
 $\theta \in \mathcal{M}$



TRANSFORMATION OF BASIS FROM DATA SPACE TO LATENT SPACE

- VAE and GAN are all based upon the approximated distance/divergence
 - VAE optimizes the ELBO
 - $\ln P(E) \geq E_{Q(H|E)} \ln P(E|H) - KL(Q(H|E) \parallel P(H))$
 - Optimization on ELBO does not guarantee the optimization of $\ln P(E)$
 - Rather, only guarantees the lower bound optimization, itself
 - Problem of approximation on $\min_{\theta \in \mathcal{M}} d(p_{data}, p_{\theta})$
 - GAN optimizes the minimax of the classification error
 - $\min_G \left\{ \max_D \left\{ E_{x \sim p_{data}(x)} [\log D(x)] + E_{z \sim p_z(z)} [\log(1 - D(G(z)))] \right\} \right\}$
 - Minimization on JS divergence is only feasible when D becomes the optimal Discriminator
 - However, it is very difficult to achieve in reality
 - Problem of approximation on $d(p_{data}, p_{\theta})$
- Then, the question becomes whether the exact inference is feasible, or not
 - Flow is a family of \mathcal{M} that enables the optimization on the exact inference
 - $\log p_{\theta}(x) = \log p_z(f_{\theta}(x)) + \log \left| \frac{\partial f_{\theta}(x)}{\partial x} \right|$



Detour: Transformation in Linear Algebra

- Imagine a linear transformation

- We can transform a basis through matrix multiplication

- $$\begin{bmatrix} x_{t+1}^1 \\ x_{t+1}^2 \end{bmatrix} = \begin{bmatrix} 2 & 1 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} x_t^1 \\ x_t^2 \end{bmatrix} \rightarrow \mathbf{x}_{t+1} = T\mathbf{x}_t$$

- The meaning of the transformation matrix

- Diagonal : scaling of the basis
- Off-diagonal : correlation of the basis
- $\text{Det}(T)$: area multiplier after the transformation
 - What-if $\text{Det}(T) < 0 \rightarrow$ The basis becomes directed inversely (flipping)

- Some matrix enables the transformation in the degenerative manner

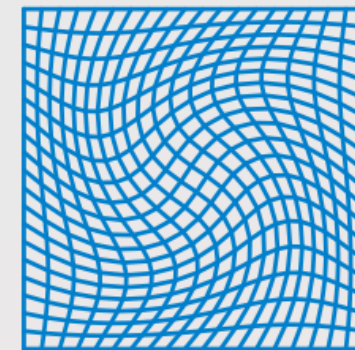
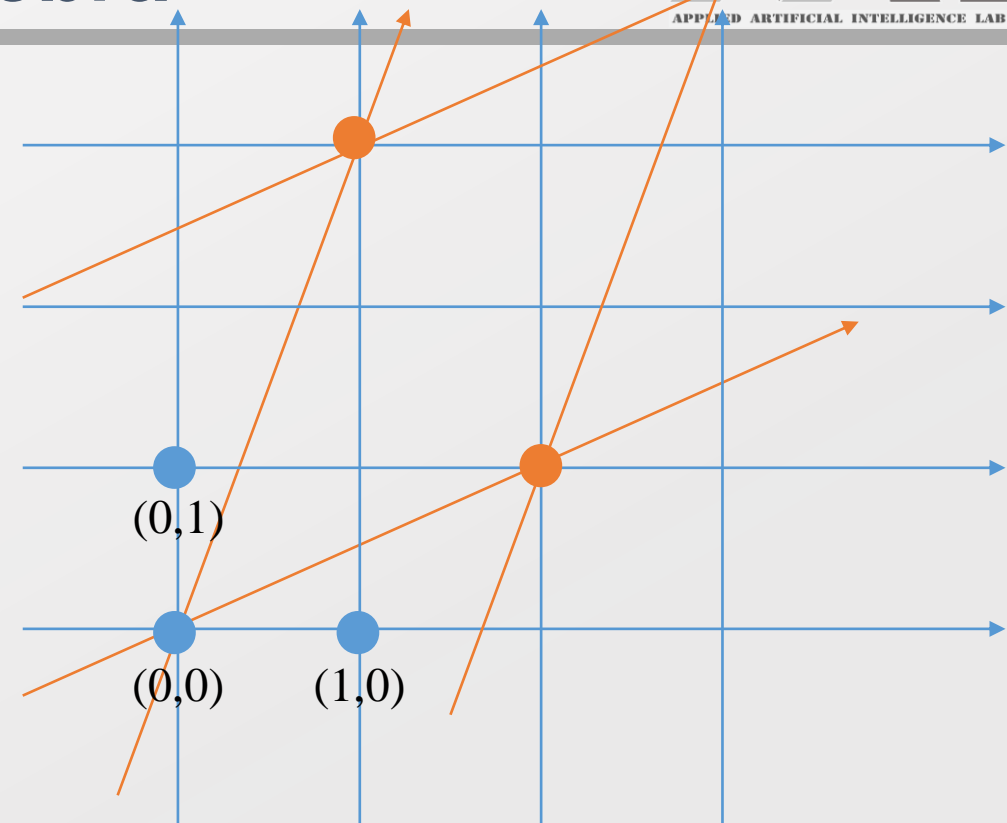
- $\text{Det}(T) = 0$
 - All points before the transformation are collapsed into a single point of (0,0)

- Linear transformation

- Can be changed as $\mathbf{x}_4 = (T_4 T_3 T_2 T_1) \mathbf{x}_0 = T_4 (T_3 (T_2 (T_1 \mathbf{x}_0))) = T_4 (T_3 (T_2 \mathbf{x}_1))$
- Can be inversed if A is not degenerative as $\mathbf{x}_0 = (T_1^{-1} T_2^{-1} T_3^{-1} T_4^{-1}) \mathbf{x}_4$

- Diffeomorphism

- Given two manifolds of M and N , a differentiable map $f: M \rightarrow N$ is a diffeomorphism
 - if f is a bijective
 - if $f^{-1}: N \rightarrow M$ is differentiable



<https://en.wikipedia.org/wiki/Diffeomorphism>

- Change of variable in a single dimension

- Condition

- Let f and h be continuous functions from $[a, b]$ to R , and assume h strictly increasing.
- Let α be the inverse function of h as $y = h(x), x = \alpha(y)$
- $a \leq x \leq b, h(a) \leq y \leq h(b)$

- Then, $\int_a^b f(x)dx = \int_{h(a)}^{h(b)} f(\alpha(y))d\alpha(y)$

- In statistics, $p_y(y) = p_x(\alpha(y))\alpha'(y)$

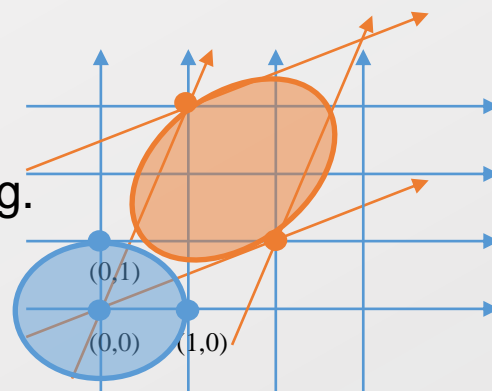
- Change of variable in its multivariate version

- $\mathbf{x}_{t+1} = T\mathbf{x}_t \rightarrow \Delta A \approx J(\mathbf{x}_t)\Delta\mathbf{x}_{t,1} \dots \Delta\mathbf{x}_{t,m}$

- \mathbf{x}_t : a m -dimensional vector
- $J(\mathbf{x}_t)$: a Jacobian matrix of T^{-1}
- A : an area in \mathbf{x}_{t+1} space

- $\iint_{X_t} f(\mathbf{x}_t)dx_{t,1} \dots dx_{t,m} = \iint_{X_{t+1}} f(T(\mathbf{x}_t))|\det J(\mathbf{x}_t)|dx_{t+1,1} \dots dx_{t+1,m}$

- In statistics, $p_{\mathbf{x}_t}(\mathbf{x}_t) = p_{\mathbf{x}_{t+1}}(T(\mathbf{x}_t))|\det J(\mathbf{x}_t)| = p_{\mathbf{x}_{t+1}}(T_\theta(\mathbf{x}_t))|\det \frac{\partial T_\theta(\mathbf{x}_t)}{\partial \mathbf{x}_t}|$

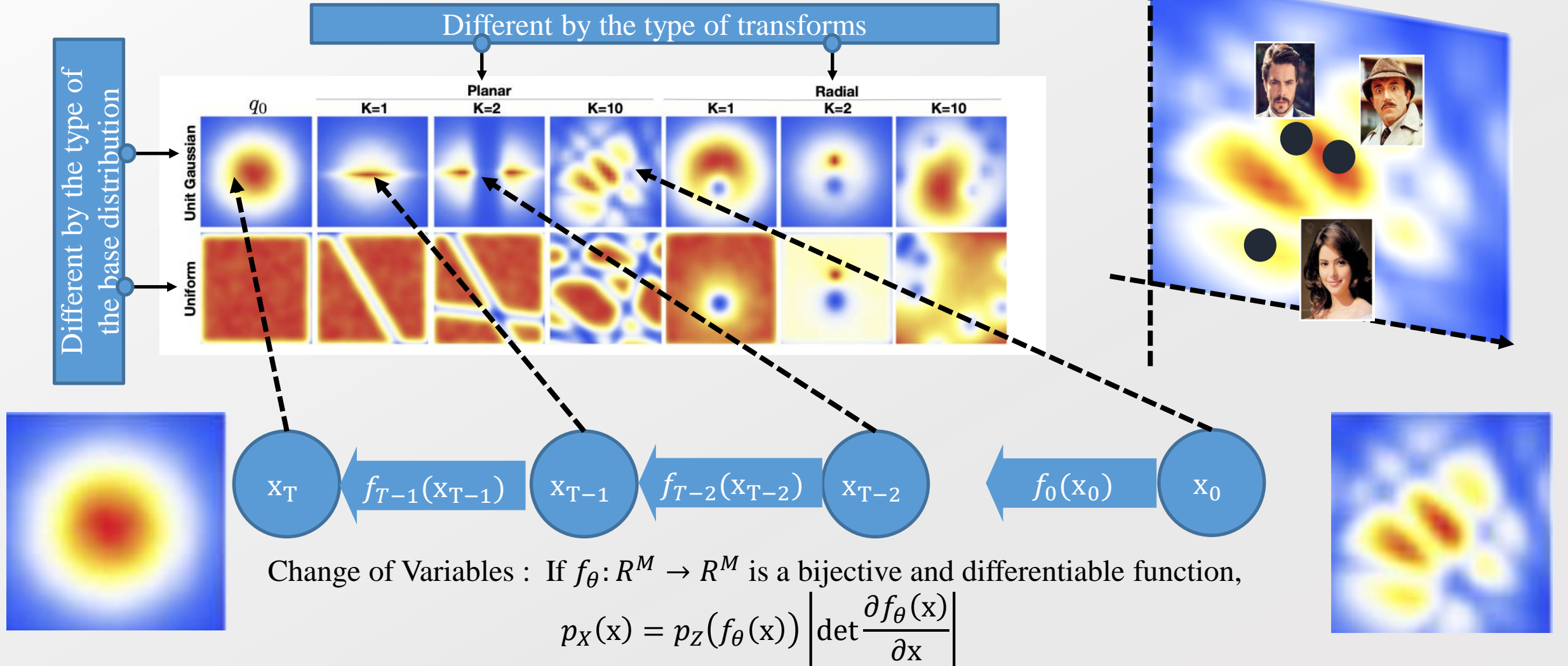


$$\text{Jacobian Martrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

assuming $f: R^n \rightarrow R^m$

$$\text{Hessian Martrix} = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 x_2} & \dots & \frac{\partial^2 f}{\partial x_1 x_n} \\ \frac{\partial^2 f}{\partial x_2 x_1} & \frac{\partial^2 f}{\partial x_2^2} & \dots & \frac{\partial^2 f}{\partial x_2 x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n x_1} & \frac{\partial^2 f}{\partial x_n x_2} & \dots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

- Generative model + Neural network
 - Neural network is used for the transformation of the distributions



Flow Model as a Formula

- Change of variable

- $x_{t+1} = T x_t, p_{x_t}(x_t) = p_{x_{t+1}}(T_\theta(x_t)) \left| \det \frac{\partial T_\theta(x_t)}{\partial x_t} \right|$

- Chain of transformations

- $x_4 = (T_4 T_3 T_2 T_1) x_0 = T_4(T_3(T_2(T_1 x_0))) = T_4(T_3(T_2 x_1))$

- Flow $f: \mathcal{X} \rightarrow \mathcal{Z}$

- An invertible differentiable function from x to z
 - Usually we denote x as data and z as noise.

- Sampling process is:

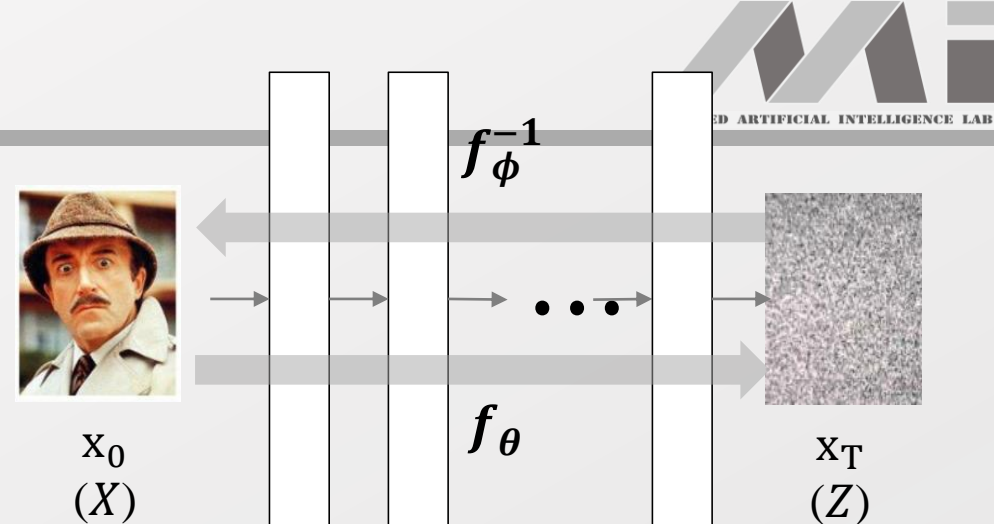
- $z \sim p(z)$
 - $x = f_\theta^{-1}(z)$

- Objective function

$$\operatorname{argmax}_\theta \mathbb{E}_{x \sim p_{data}} [\log p_\theta(x)]$$

- How can we evaluate the marginal likelihood $p_\theta(x)$? \rightarrow Using change-of-variables

$$\log p_X(x) = \log p_Z(f_\theta(x)) + \log \left| \frac{\partial f_\theta(x)}{\partial x} \right|$$



- Each transformation needs to be diffeomorphism
 - Difficult to create highly nonlinear behavior with a simple transformation
- Compose the multiple transformation as a chain
 - $f = f_K \circ f_{K-1} \circ \dots \circ f_1$
 - Sampling can be $x = f^{-1}(z) = (f_1^{-1} \circ \dots \circ f_{K-1}^{-1} \circ f_K^{-1})(z)$
 - Density evaluation
 - $p_X(x) = p_Z(f_\theta(x)) \left| \det \frac{\partial f_\theta(x)}{\partial x} \right|$
 - $\det \frac{\partial f(x)}{\partial x} = \det \prod_{i=1}^K \frac{\partial f_i}{\partial f_{i-1}} = \prod_{i=1}^K \det \frac{\partial f_i}{\partial f_{i-1}}$
 - Training with log-likelihood

$$\max_{\theta} \left\{ \log p_Z(f_\theta(x)) + \sum_{i=1}^K \log \left| \det \frac{\partial f_i}{\partial f_{i-1}} \right| \right\}$$

- Then, our question becomes what could be the choice of f
 - Following the diffeomorphism
 - Generating a flexible and domain-adapted transformation, i.e. the locality on image

STRUCTURE OF NORMALIZING FLOW

- Elementwise Flows

$$f(x) = f((x_1, \dots, x_D)) = (f_1(x_1), \dots, f_D(x_D)) = z$$

- Each dimension takes only the same dimension value as the transformation input
 - x_D : D -th feature dimension of an input vector of x
 - $z_d = f_d(x_d) = f(x_d; \theta_d)$
 - θ_d : Each dimension's transformation is parameterized per each dimension with the shared structure
 - We will treat z as the transformed x

- f is invertible if all $f^{(i)}$ be invertible.

- Analytic Inverse : $f^{-1}((z_1, \dots, z_D)) = (f_1^{-1}(z_1), \dots, f_D^{-1}(z_D))$

- Jacobian matrix

- $J_f(x) = \text{diag}(f'_1(x_1), \dots, f'_D(x_D)) \rightarrow \det(J_f(x)) = \prod_{i=1}^D f'_i(x_i)$

- Given the objective function of $\max_{\theta} \left\{ \log p_Z(f_{\theta}(x)) + \sum_{i=1}^K \log \left| \det \frac{\partial f_i}{\partial f_{i-1}} \right| \right\}$

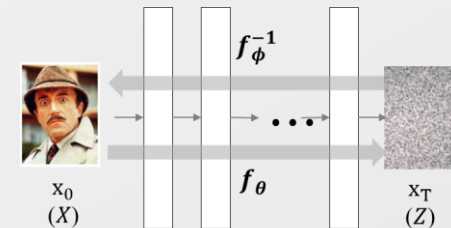
- No mixing of variables

- i.e. images have highly localized correlations \rightarrow Problem in modeling the density of the joint space

- But from these, we know that it is okay to use some activation functions or batch normalization in composition of flows.

- Computation cost of Inverse: $O(D)$
 - Computational cost of determinant: $O(D)$

- Equivalent to only applying the activation functions without any cross links between inputs and outputs



$$Jacobian\ Matrix = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

- Linear Flows (or Affine Flows)

$$f(x) = f((x_1, \dots, x_D)) = Ax + b = z$$

$$\text{Jacobian Martrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

- f is invertible if A is invertible.
 - Analytic inverse : $f^{-1}(z) = A^{-1}(z - b)$
 - Jacobian matrix
 - $J_f(x) = A \rightarrow \det(J_f(x)) = \det A$
- Given the objective function of $\max_{\theta} \left\{ \log p_Z(f_{\theta}(x)) + \sum_{i=1}^K \log \left| \det \frac{\partial f_i}{\partial f_{i-1}} \right| \right\}$
 - No capability of nonlinear transformation
 - i.e. Let's assume $z \sim N(0, I)$, then $x \sim N(\mu, \Sigma)$: This nature will not change by composing the transformation
 - Why? This is a linear transformation without any activation
 - If we add an activation to this transformation, the inversion and the Jacobian will not be simple
 - High cost of calculation
 - Computation cost of Inverse: $O(D^3)$
 - Computational cost of determinant: $O(D^3)$
 - Equivalent to applying the fully connected neural network layer without activation

- The objective function of flow
 - $\max_{\theta} \left\{ \log p_Z(f_{\theta}(x)) + \sum_{i=1}^K \log \left| \det \frac{\partial f_i}{\partial f_{i-1}} \right| \right\}$
 - The calculation of the determinant of the Jacobian matrix becomes the key
- Laplace expansion is an expression of determinant as the sum of small matrixes

- Given the $n \times n$ matrix of A
 - Let's say $M_{i,j}$ is the determinant of the minor matrix which removes i and j -th row and column of A
 - Then, the determinant of A becomes

$$\det(A) = \sum_{j=1}^n (-1)^{i+j} A_{i,j} M_{i,j}$$

- This suggests that the determinant calculation becomes simpler with a well positioned $A_{i,j} = 0$
 - Diagonal matrix and Triangle matrix

$$\det(A) = \prod_{i=1}^n A_{i,i}$$

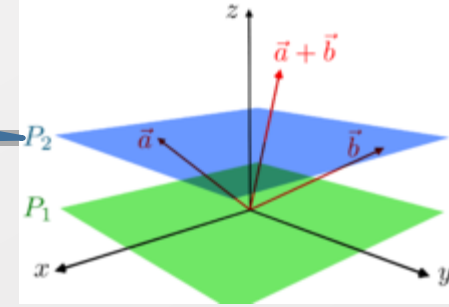
- Computation complexity
 - Laplace expansion takes $O(n!)$
 - LU decomposition takes $O(n^3)$

$$\text{Jacobian Martrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

$$\text{Triangular Jacobian Martrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} = 0 & \dots & \frac{\partial f_1}{\partial x_n} = 0 \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & & \frac{\partial f_2}{\partial x_n} = 0 \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \dots & \frac{\partial f_m}{\partial x_n} = 0 \end{bmatrix}$$

Coupling Flows: Affine

Affine subspace



- Only transform the affine subspace of x

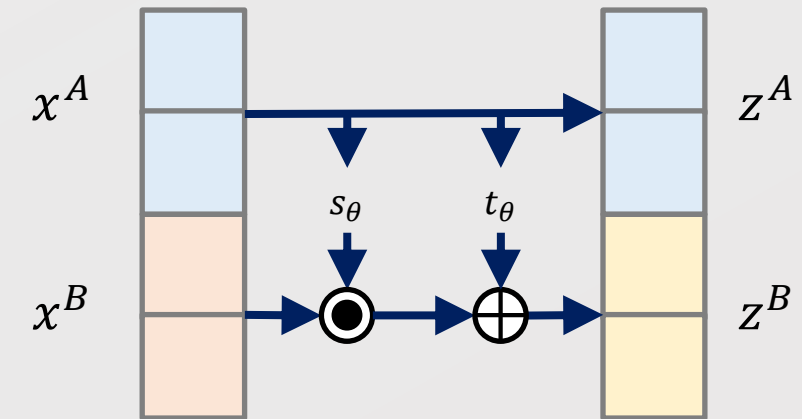
Affine Coupling Flows

- Split variables into x^A and $x^B \rightarrow$ elementwise flows on x^B with the effect of x^A

$$f(x) = f((x_1, \dots, x_D)) = f([x^A; x^B]) = (x^A, x^B \cdot s_\theta(x^A) + t_\theta(x^A)) = z$$

- In other words, $f = (f^A, f^B)$ where $z^A = f^A(x^A) = x^A$ and $z^B = f^B(x^B) = x^B \cdot s_\theta(x^A) + t_\theta(x^A)$
 - $z^A = x^A$: No transformation
 - $z^B = x^B \cdot s_\theta(x^A) + t_\theta(x^A)$
- Utilizing the elementwise flow in x^B : Good characteristics in nonlinearity and computation complexity
- f is invertible without significant restrictions of s_θ and t_θ
 - Analytic inverse
 - $x^A = (f^A)^{-1}(z^A) = z^A$
 - $x^B = (f^B)^{-1}(z^B) = \{z^B - t_\theta(z^A)\} / s_\theta(z^A)$
- Jacobian Matrix

$$J_f(x) = \begin{bmatrix} \frac{\partial f^A}{\partial x^A} & \frac{\partial f^A}{\partial x^B} \\ \frac{\partial f^B}{\partial x^A} & \frac{\partial f^B}{\partial x^B} \end{bmatrix} = \begin{bmatrix} I & 0 \\ \frac{\partial f^B}{\partial x^A} & \text{diag}(s(x^A)) \end{bmatrix} \Rightarrow \det J_f(x) = \prod_j s(x^A)_j$$



- Affine Coupling Flows

- Split variables into x^A and $x^B \rightarrow$ elementwise flows on x^B with the effect of x^A

$$f(x) = f((x_1, \dots, x_D)) = f([x^A; x^B]) = (x^A, x^B \cdot s_\theta(x^A) + t_\theta(x^A)) = z$$

- Given the objective function of $\max_{\theta} \left\{ \log p_Z(f_\theta(x)) + \sum_{i=1}^K \log \left| \det \frac{\partial f_i}{\partial f_{i-1}} \right| \right\}$

- Some capability of nonlinear transformation

- $s_\theta(x^A)$ and $t_\theta(x^A)$: the source of the nonlinear nature
- i.e. Let's assume $z \sim N(0, I)$, then $x \sim N(\mu, \Sigma)$: This nature will not change by composing the transformation
- Why? This is a linear transformation without any activation
 - If we add an activation to this transformation, the inversion and the Jacobian will not be simple

- Low cost of calculation

- Computation cost of Inverse: $O(D)$
- Computational cost of determinant: $O(D)$

- Equivalent to applying the some selective links on neural network layer with some activations

- Moreover, through the compositions of transformations, the split can be changed

- Split variables into x^A and x^B at $f_1 \neq$ Split variables into x^A and x^B at f_2

- Affine Coupling Flows

- Split variables into x^A and $x^B \rightarrow$ elementwise flows on x^B with the effect of x^A

$$f(x) = f((x_1, \dots, x_D)) = f([x^A; x^B]) = (x^A, x^B \cdot s_\theta(x^A) + t_\theta(x^A)) = z$$

- Generalized Coupling Flows : general approach to construct non-linear flows

$$f(x) = f([x^A; x^B]) = (x^A, g(x^B; \theta(x^A))) = z$$

- i.e. $f = (f^A, f^B)$ where $f^A(x^A) = x^A$ (identity), $f^B(x^B) = g(x^B; \theta(x^A))$

- $\theta(x^A)$: the parameterization dynamically determined by x^A

- f is invertible if g is invertible. (not θ !)

- Analytic inverse

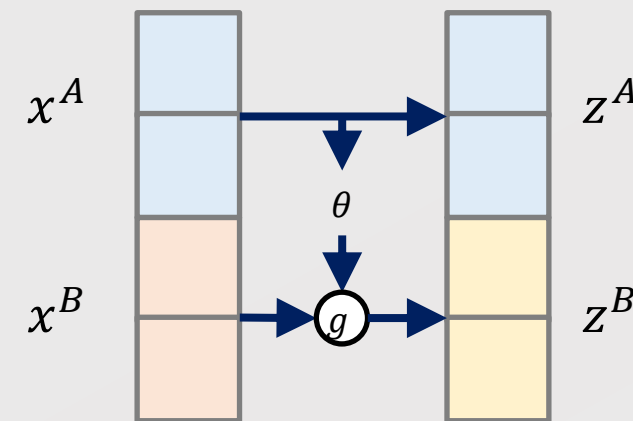
- $f^{-1}(z) = (z^A, g^{-1}(z^B; \theta(z^A)))$

- Jacobian Matrix

- $J_f(x) = \begin{bmatrix} \frac{\partial f^A}{\partial x^A} & \frac{\partial f^A}{\partial x^B} \\ \frac{\partial f^B}{\partial x^A} & \frac{\partial f^B}{\partial x^B} \end{bmatrix} = \begin{bmatrix} I & 0 \\ \frac{\partial g}{\partial x^A} & J_g(x^B) \end{bmatrix} \Rightarrow \det J_f(x) = \det J_g(x^B) = \prod_j J_g(x^B)_j$

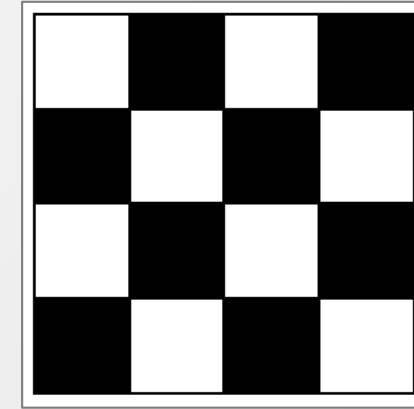
- g can be any function to be applied in the elementwise manner

- Linear function, Spline function, and inverse CDF....

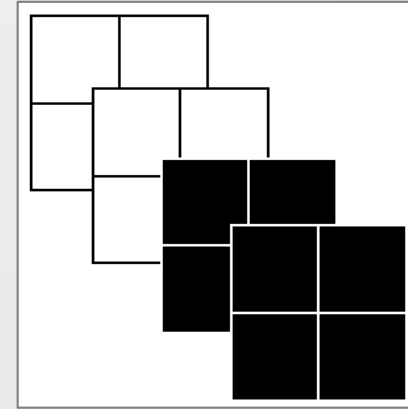


- Coupling flow

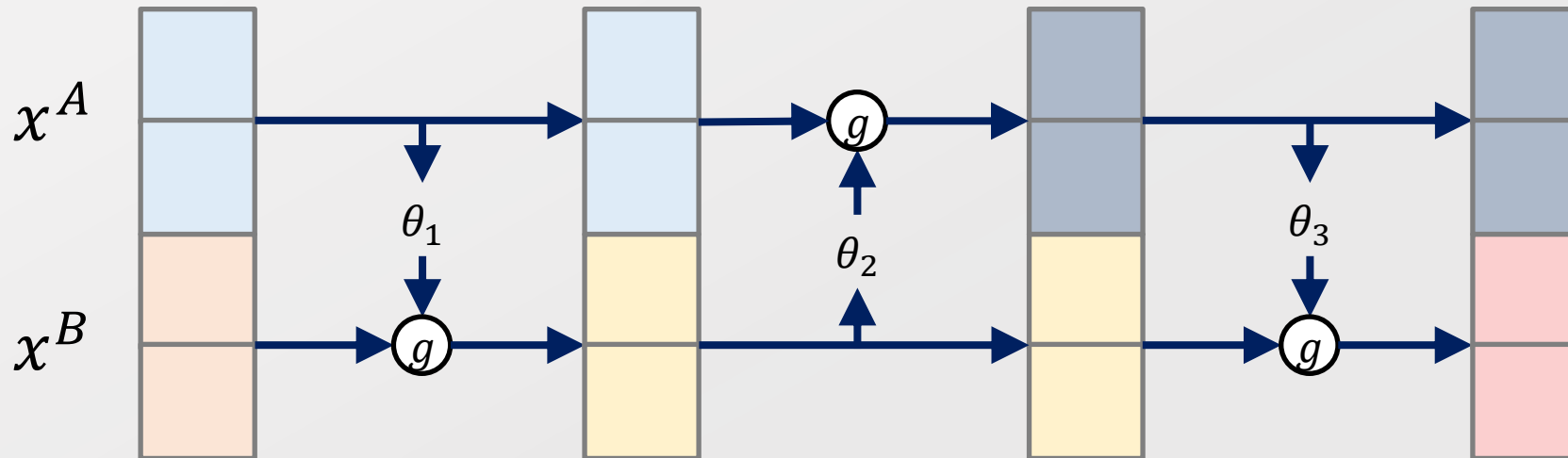
- Split variables into x^A and $x^B \rightarrow$ elementwise flows on x^B with the effect of x^A
- Setting x^A and x^B depends upon the domain of application
 - Images? Sequences?....
- Split needs to change over the composition of the transformation



Spatial Checkerboard Pattern



Channel-wise Masking



- Coupling Flows (a.k.a. RealNVP) : general approach to construct non-linear flows

$$f(x) = f([x^A; x^B]) = (x^A, g(x^B; \theta(x^A))) = z$$

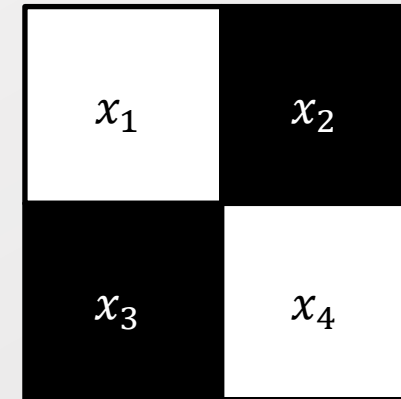
- Interpretation from the probabilistic modeling perspective
 - Flow of each variable becomes the Bayesian network links

- $f(x_2) = x_2$
- $f(x_3) = x_3$
- $f(x_1) = x_1 s_1(x_2, x_3) + t_1(x_2, x_3)$
- $f(x_4) = x_4 s_4(x_2, x_3) + t_4(x_2, x_3)$

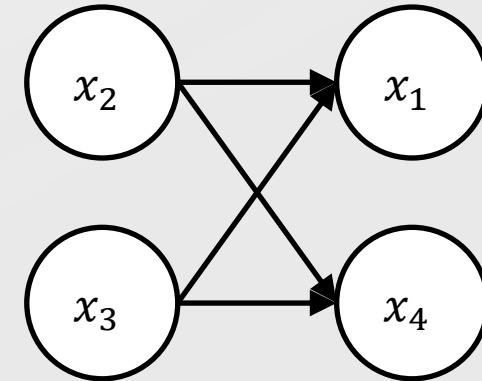
- Is same as

- $f(x_2) = x_2$
- $f(x_3) = x_3$
- $f(x_1) = x_1 s_1(\text{parent}(x_1)) + t_1(\text{parent}(x_1))$
- $f(x_4) = x_4 s_4(\text{parent}(x_4)) + t_4(\text{parent}(x_4))$

- This creates a DAG, which is the key characteristics of Bayesian network
 - However, very typical and stereotype
 - Then, is it possible to create a more interesting structure by Bayesian network?
 - with some inspiration from the old graphical models?

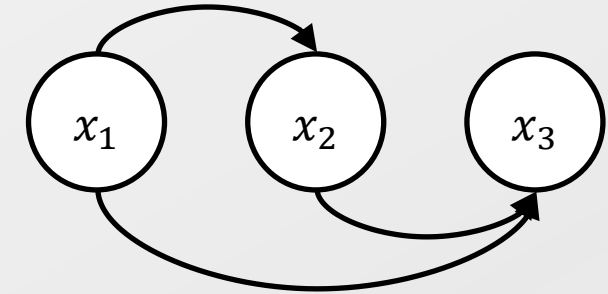


Variable Split



Bayesian Network

- Autoregressive Bayesian network structure
 - Density model
 - $x_1 \sim p_\theta(x_1)$
 - $x_2 \sim p_\theta(x_2|x_1)$
 - $x_3 \sim p_\theta(x_3|x_1, x_2)$
 - Generation process : sampling x given z : $p(x|z)$
 - $x_1 = f_\theta^{-1}(z_1)$
 - $x_2 = f_\theta^{-1}(z_2; x_1)$
 - $x_3 = f_\theta^{-1}(z_3; x_1, x_2)$
 - Inference process : learning the distribution of z given x : $q(z|x)$
 - $z_1 = f_\theta(x_1)$
 - $z_2 = f_\theta(x_2; x_1)$
 - $z_3 = f_\theta(x_3; x_1, x_2)$
- Jacobian matrix : triangular matrix because of independence between z_t and $x_{(t+1):T}$
- Computation complexity
 - Let's consider the order
 - Inference process : all inputs are given : Can be parallel \rightarrow Fast
 - Generation process : some inputs are given : Needs to be sequential \rightarrow Slow



- Autoregressive flow

- Generation process : sampling x given $z : p(x|z)$
 - $x_1 = f_{\theta}^{-1}(z_1), x_2 = f_{\theta}^{-1}(z_2; x_1), x_3 = f_{\theta}^{-1}(z_3; x_1, x_2)$
- Inference process : learning the distribution of z given $x : q(z|x)$
 - $z_1 = f_{\theta}(x_1), z_2 = f_{\theta}(x_2; x_1), z_3 = f_{\theta}(x_3; x_1, x_2)$

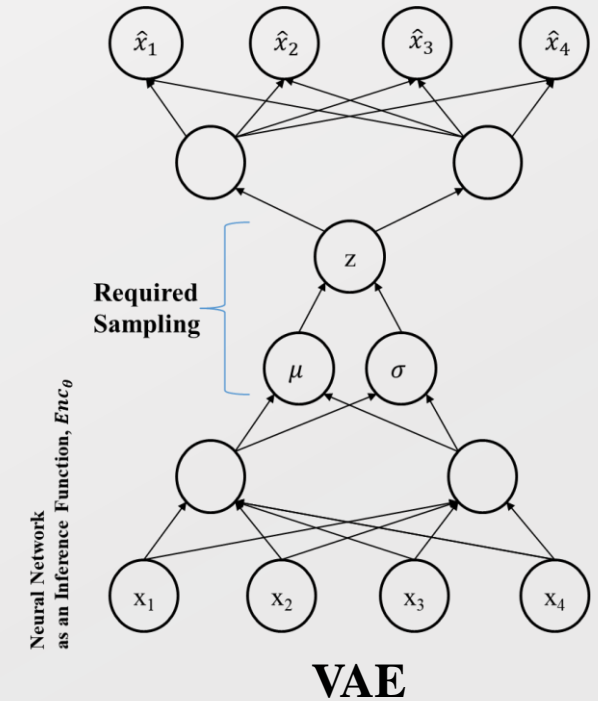
- Inverse autoregressive flow

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

- Interpretation by the latent inference algorithm, such as variational inference

- $$\log p(x) \geq \mathbb{E}_{z \sim q(z|x)} [\log p(x|z)] - D_{KL}[q(z|x) || p(z)] = \mathbb{E}_{z \sim q(z|x)} [\log p(x|z)] - \mathbb{E}_{z \sim q(z|x)} \left[\frac{q(z|x)}{p(z)} \right]$$

$$\approx E_Z \left[\log p(x|z) - \left(\frac{q(z|x)}{p(z)} \right) \right] \text{ where } z \sim q(z|x)$$
 - Finding $z \sim q(z|x) \rightarrow$ Flexible Posterior by $z_1 = f_{\theta}(x_1), z_2 = f_{\theta}(x_2; z_1), z_3 = f_{\theta}(x_3; z_1, z_2)$
- In VAE, Need the generation on $z \sim q(z|x)$ to calculate the expectation in the Monte-carlo manner



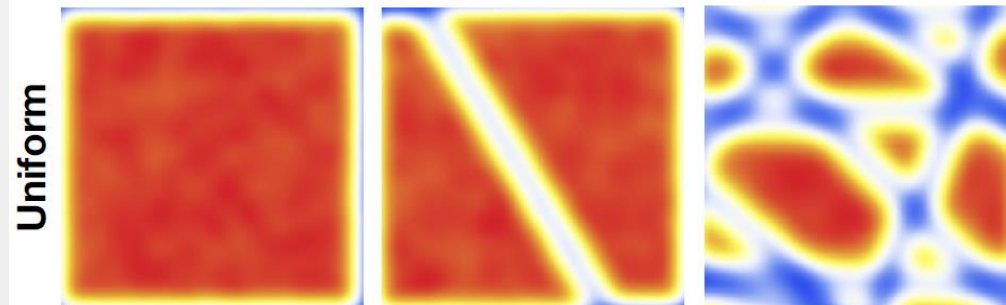
- What could be the structure of the transformation function?
 - By following the Bayesian network interpretation
 - $x_i = f_{\theta}(z_i; \text{parent}(x_i))$
- Desired/Required nature of f_{θ}
 - Need to be applied to a single dimension by creating influence from the other many
 - Invertible
 - Nonlinear
 - Variations
 - NICE, RealNVP, IAF-VAE
 - Using affine transformation (common)
 - Glow
 - Using invertible 1x1 convolution
 - Flow++
 - Using continuous mixture CDFs
 - Spline Flow
 - Using splines

Diverse Choices of Transformation

	Name	$f(x)$	(f^{-1}, Df) Complexity
Finite Flow	Elementwise Flow	$(f(x_1), \dots, f(x_D))$	
	Linear Flow	$Ax + b$	$O(D^3), O(D^3)$
	Planar Flow	$x + uh(w^T x + b)$	$-, O(D)$
	Radial Flow	$x + \frac{\beta}{\alpha + \ x - x_0\ } (x - x_0)$	$-, O(D)$
	Coupling Flow	$(h(x^A; \Theta(x^B)), x^B)$	Invertible iff h is invertible / Depends on h , but efficient
	Autoregressive Flow	$(h(x_t; \Theta_t(x_{1:t-1})))_t$	Inverse cost high / determinant elementwise product
	Spline Flow	Spline	Inverse requires Newton method / determinant closed form
	Residual Flow	$(x^A + F(x^B), x^B + G(f(x)^A))$	Invertible easy / Jacobian difficult
Infinitesimal Flow	Langevin Flow	$\frac{\partial}{\partial t} q_t(x) = - \sum_i \frac{\partial}{\partial x_i} [F_i(x, t) q_t] + \frac{1}{2} \sum \frac{\partial}{\partial x_i \partial x_j} [D_{ij}(x, t) q_t]$	
	Hamiltonian Flow	$\mathcal{H}(x, \omega) = -\mathcal{L}(x) - \frac{1}{2} \omega^T M \omega$	

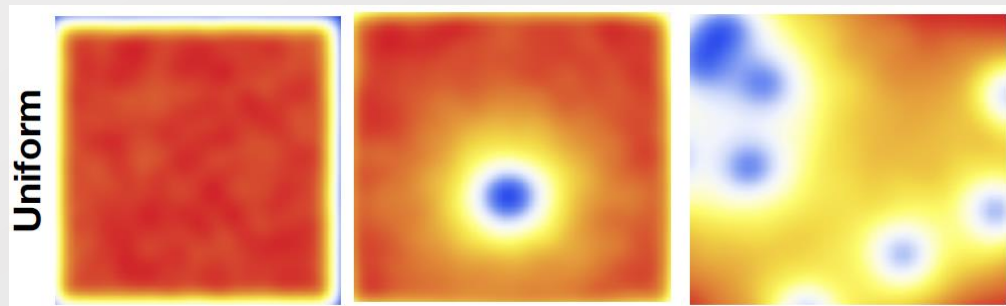
- Planar Flow

- $f_k(z_{k-1}) = z_{k-1} + u_k h_k(w_k^T z_{k-1} + b_k)$
 - z_{k-1} : Linear component
 - $u_k h_k(w_k^T z_{k-1} + b_k)$: Nonlinear component
- Determinant of Jacobian
 - $\left| \frac{\partial f_k}{\partial z_{k-1}} \right| = \left| I + u_k (h'_k(w_k^T z_{k-1} + b_k) w_k)^T \right| = \left| 1 + u_k^T h'_k(w_k^T z_{k-1} + b_k) w_k \right|$
- Objective function
 - $\log q_K(z_K) = \log q_0(z_0) - \sum_{k=1}^K \log |1 + u_k^T \psi_k(z_{k-1})|$, where $\psi_k(z_{k-1}) = h'_k(w_k^T z_{k-1} + b_k) w_k$



- Radial Flow

- $f_k(z_{k-1}) = z_{k-1} + \frac{\beta}{\alpha + \|z_{k-1} - z_{k-1}^*\|} (z_{k-1} - z_{k-1}^*)$
- Determinant of Jacobian
 - $\left| \frac{\partial f_k}{\partial z_{k-1}} \right| = \left[1 + \frac{\beta}{\alpha + \|z_{k-1} - z_{k-1}^*\|} \right]^{D-1} \left[1 + \frac{\beta}{\alpha + \|z_{k-1} - z_{k-1}^*\|} - \frac{\|z_{k-1} - z_{k-1}^*\|}{(\alpha + \|z_{k-1} - z_{k-1}^*\|)^2} \right]$
- Objective function
 - $\log q_K(z_K) = \log q_0(z_0) - \sum_{k=1}^K (D-1) \log \left(1 + \frac{\beta}{\alpha + \|z_{k-1} - z_{k-1}^*\|} \right) - \sum_{k=1}^K \log \left(1 + \frac{\beta}{\alpha + \|z_{k-1} - z_{k-1}^*\|} - \frac{\|z_{k-1} - z_{k-1}^*\|}{(\alpha + \|z_{k-1} - z_{k-1}^*\|)^2} \right)$



COMBINATION WITH OTHER GENERATIVE MODELS

- Rank-Nullity Theorem

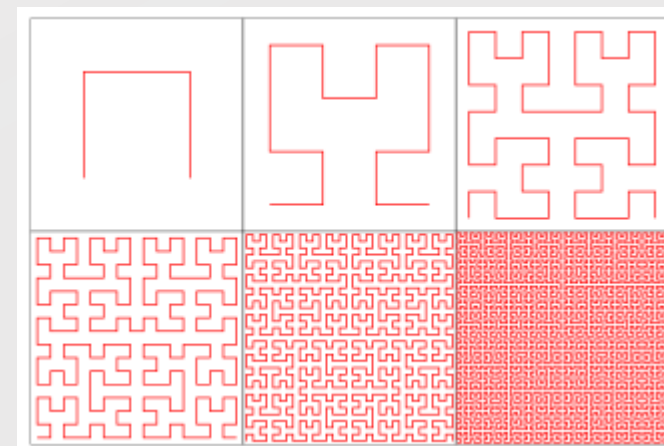
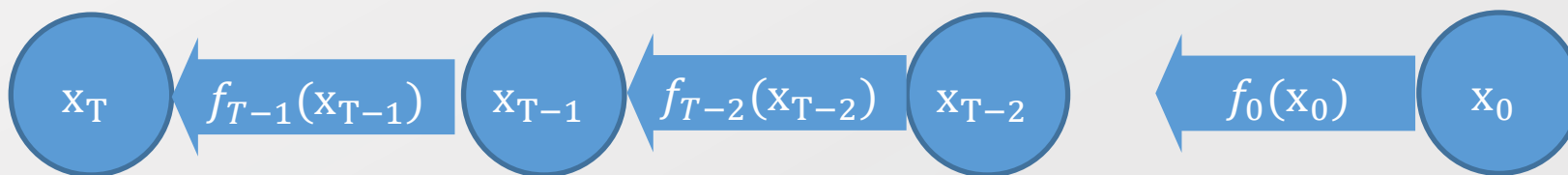
- [Rank-Nullity Theorem] Let $T: V \rightarrow W$ be a linear transformation, where V and W are the finite-dimensional vector spaces. Then, $\text{Rank}(T) + \text{Nullity}(T) = \dim V$, where $\text{Rank}(T) = \dim(T(V))$ and $\text{Nullity}(T) = \dim(\text{Ker}(T))$.

- Simple interpretations

- [Theorem] A function $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ with $n > m$ **can** be bijective.
- [Theorem] A **continuous** function $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ with $n > m$ **cannot** be bijective.

- Requirements of the flow transformation

- Bijjective function
- Differentiable function
- Hence, the flow transformation needs to be the matching dimension between the data and the latent variables



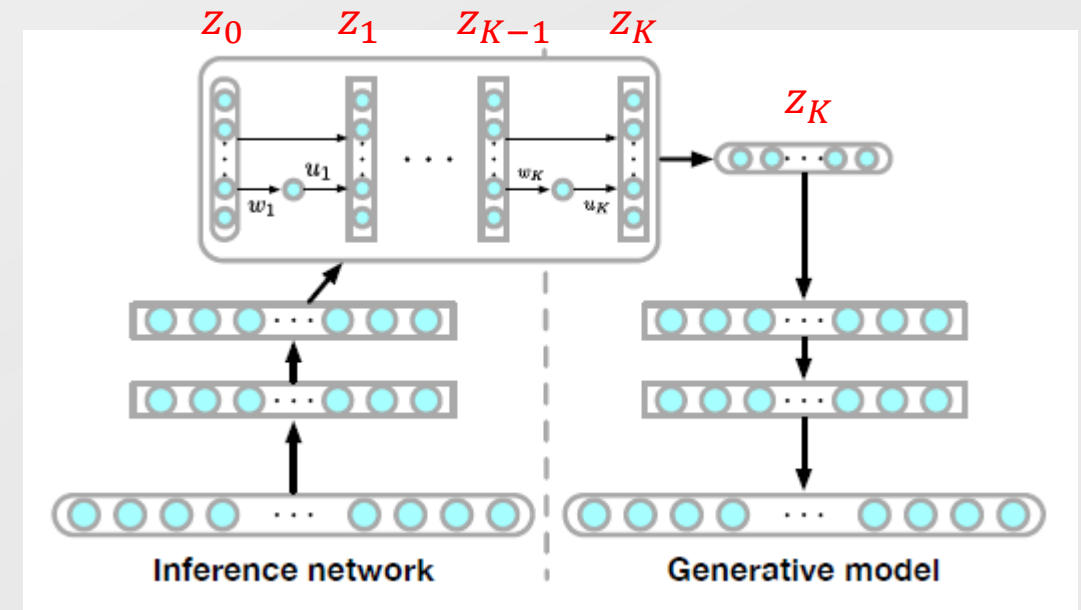
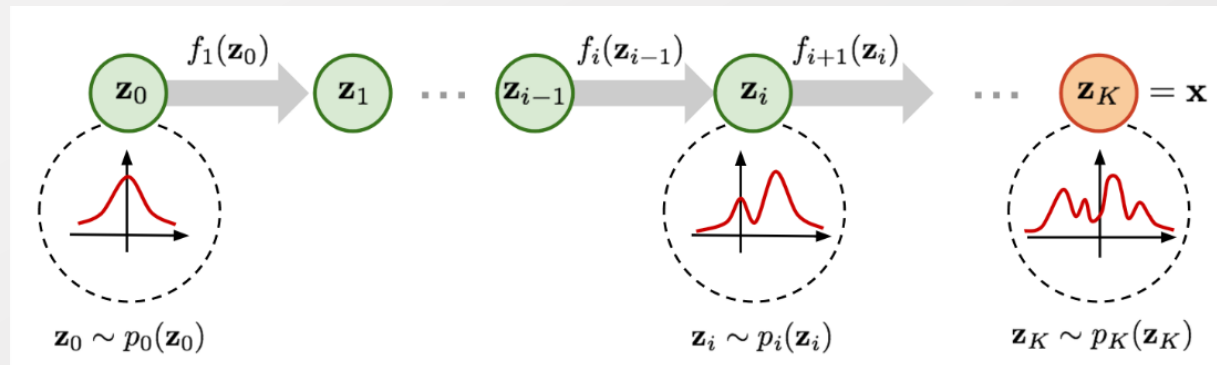
Hilbert Curve

(<https://math.stackexchange.com/questions/557801/why-is-this-not-a-space-filling-curve>)

- $\mathbb{E}_{q_{\phi}(h|e)} \left[\ln \frac{q_{\phi}(h)}{p(h)} \right]$ is the expected density ratio
 - The density ratio can be computed by building a classifier to distinguish observed data from that generated by the model.
 - $p(h)$ is used for a set of n samples $E_p = \{e_1^{(p)}, \dots, e_n^{(p)}\}$
 - $q_{\phi}(h)$ is used for a set of n samples $E_q = \{e_1^{(q)}, \dots, e_n^{(q)}\}$
 - Through the ancestral sampling of $q_{\phi}(h|e)$ by randomly selecting e
 - A random variable y that assigns a label $y = 1$ to all samples in E_p and $y = 0$ to all samples in E_q .
 - $p(h) = p(h|y = 0)$ and $q_{\phi}(h) = p(h|y = 1)$
 - $p^*(h|y) \equiv \begin{cases} q_{\phi}(h), y = 1 \\ p(h), y = 0 \end{cases}$
- By applying Bayes' rule, we can compute the ratio $r(\mathbf{h}) = \ln \frac{q_{\phi}(h)}{p(h)}$ as:
 - $$\frac{q_{\phi}(h)}{p(h)} = \frac{p^*(h|y=1)}{p^*(h|y=0)} = \frac{\frac{p^*(h,y=1)}{p^*(y=1)}}{\frac{p^*(h,y=0)}{p^*(y=0)}} = \frac{\frac{p^*(y=1|h)p^*(h)}{p^*(y=1)}}{\frac{p^*(y=0|h)p^*(h)}{p^*(y=0)}} = \frac{p^*(y=1|h)}{p^*(y=0|h)} \cdot \frac{\pi}{1-\pi} = \frac{p^*(y=1|\mathbf{h})}{p^*(y=0|\mathbf{h})} = \frac{D(h)}{1-D(h)}$$
 - $D(h) = p^*(y = 1|h)$
 - Which indicates density ratio estimation = class probability estimation.
 - The problem is reduced to computing the probability $p(y = 1|h)$
 - Discriminative modeling can be applied, as a discriminator

Discriminator in VAE
for
Optimal Prior Learning

- Merge of two generative models
 - Variational autoencoder
 - requires a flexible latent space model
 - provides the information extraction through data compressions
 - Flow
 - restricted to the continuous and bijective transformation, so no dimensionality reduction
 - provides the flexible modeling on the probability space
- What-if we model the latent space by flow
 - If $z_0 \sim q_0(z_0)$ and
 - $z_K = f_K \circ \dots \circ f_2 \circ f_1(z_0)$,



- Original ELBO in VAE

- $\mathcal{L} = \mathbb{E}_{q_\phi(H|E)}[\log p_\theta(E|H)] - D_{KL}(q_\phi(H|E) || p_\theta(H))$

- $$L_{ELBO}(\theta, \phi) = \int q_\phi(z_K|x) \log p_\theta(x|z_K) dz_K + \int q_\phi(z_K|x) \log \frac{p(z_K)}{q_\phi(z_K|x)} dz_K$$

$$= E_{z_K \sim q_\phi(z_K|x)} [\log p_\theta(x|z_K)] + E_{z_K \sim q_\phi(z_K|x)} \left[\log \frac{p(z_K)}{q_\phi(z_K|x)} \right]$$

$$= E_{z_0 \sim q_\phi(z_0|x)} [\log p_\theta(x|f_K \circ \dots \circ f_1(z_0))] + E_{z_0 \sim q_\phi(z_0|x)} \left[\log \frac{p(z_K)}{q_\phi(f_K \circ \dots \circ f_1(z_0)|x)} \right]$$

$$= E_{z_0 \sim q_\phi(z_0|x)} [\log p_\theta(x|f_K \circ \dots \circ f_1(z_0))] + E_{z_0 \sim q_\phi(z_0|x)} [\log p(z_K)]$$

$$- E_{z_0 \sim q_\phi(z_0|x)} [\log q_\phi(f_K \circ \dots \circ f_1(z_0)|x)]$$

- Law of the unconscious statistician (LOTUS)

- $E_{z_K \sim q_K(z_K)} [\log h(z_K)] = E_{z_0 \sim q_0(z_0|x)} [\log h(f_K \circ \dots \circ f_1(z_0))]$

- Flow

- $\log q_\phi(z_K) = \log q_\phi(z_0) - \sum_{k=1}^K \log \det |D_{z_k} f_k|$

- Eventually, ELBO becomes

- $$L_{ELBO}(\theta, \phi) = E_{z_0 \sim q_\phi(z_0|x)} [\log p_\phi(x|f_K \circ \dots \circ f_1(z_0))] + E_{z_0 \sim q_\phi(z_0|x)} [\log p(z_K)]$$

$$- E_{z_0 \sim q_\phi(z_0|x)} [\log q_\phi(z_0|x)] + \sum_{k=1}^K E_{z_0 \sim q_\phi(z_0|x)} [\log \det |D_{z_k} f_k|]$$

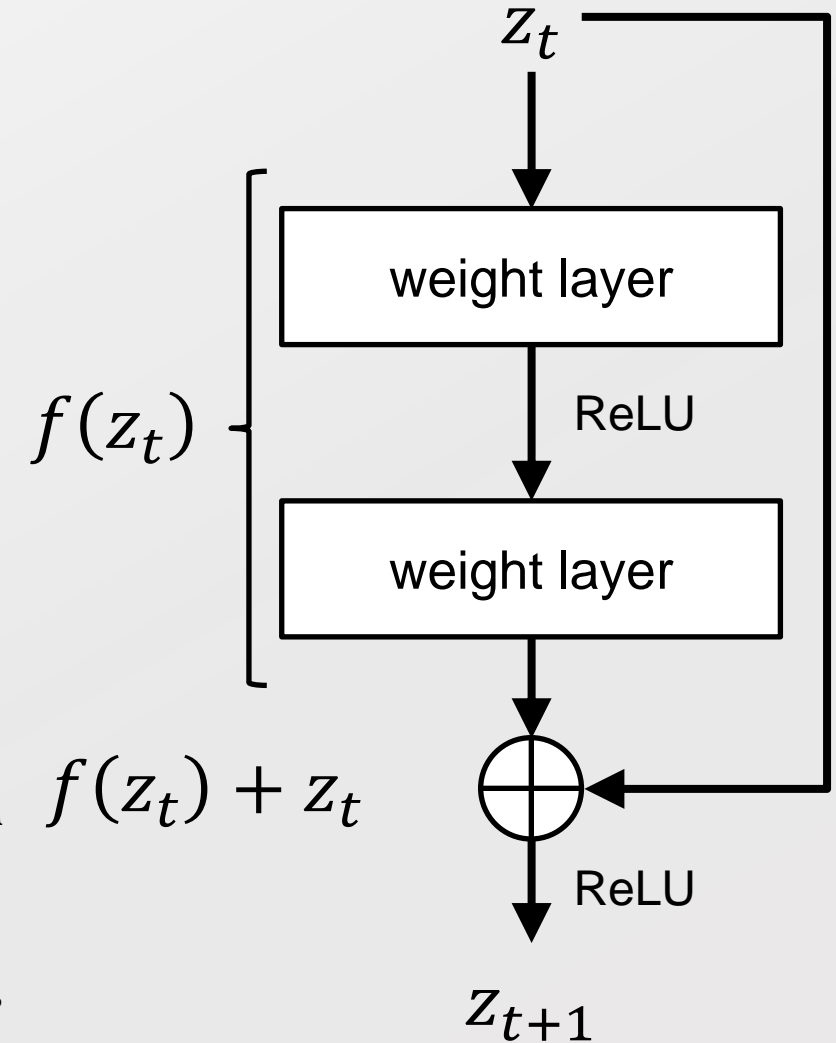
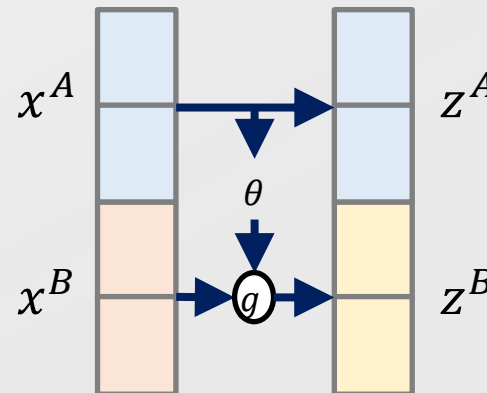
CONTINUOUS TRANSFORMATION ON FLOW

- Generalized Coupling Flows : general approach to construct non-linear flows

$$f(x) = f([x^A; x^B]) = (x^A, g(x^B; \theta(x^A))) = z$$

- $f = (f^A, f^B)$
 - where $f^A(x^A) = x^A$ (identity), $f^B(x^B) = g(x^B; \theta(x^A))$
- Residual Network Structure
 - $F(z_t) = W_2 \sigma(W_1 z_t)$
 - $z_{t+1} = z_t + f(z_t; \{W_i\})$
- Similarity
 - Existence of the identity operation

- ResNet
 - All dimensions are processed
- Flow
 - Some dimensions are processed
 - The other dimensions are kept



- Invertible Residual Networks

- Generative Modeling with i-ResNets

- Let $F: \mathcal{X} \rightarrow \mathcal{Z}$ be a transformation function. Then, equation of flow models:

$$\log p_x(x) = \log p_z(z) + \log \left| \det \frac{\partial F}{\partial x} \right|$$

- From the residual connection ($F = I + f$),

$$\log \left| \det \frac{\partial F}{\partial x} \right| = \text{tr} \left(\log \left(I + \frac{\partial f}{\partial x} \right) \right) = \sum_{k=1}^{\infty} (-1)^{k+1} \frac{\text{tr} \left(\left(\frac{\partial f}{\partial x} \right)^k \right)}{k}$$

- Stochastic Approximation of log-determinant

- 1) Hutchinson's trace estimator $\rightarrow \text{tr}(J_f) = \mathbb{E}_{p(v)}[v^T J_f v]$ with $\mathbb{E}[v] = 0, \text{Cov}(v) = I$
- 2) Infinite sum \rightarrow truncated at index n

$$\log p_x(x) = \log p_z(z) + \mathbb{E}_v \left[\sum_{k=1}^n \frac{(-1)^{k+1}}{k} v^T J_f(x)^k v \right]$$

- By 2), the estimator is biased. (but the paper said 5-10 iterations are enough)

- Ordinary Differential Equation

- Differential equation with a set of functions on an independent variable and the derivatives of the functions

- $F(x, y, y', \dots, y^{(n)}) = 0$

- y : dependent variable
- x : independent variable
- $y = f(x)$: unknown function of x

- System of Differential Equations

- Coupled differential equations

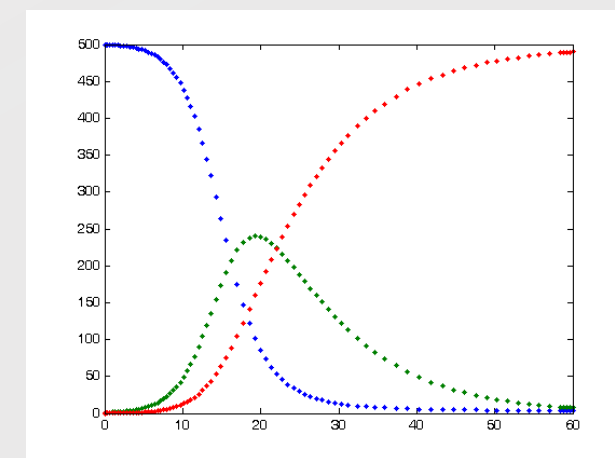
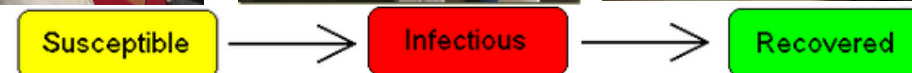
- $\mathbf{y} : \mathbf{y}(x) = [y_1(x), y_2(x), \dots, y_m(x)]$

- $$\begin{pmatrix} f_1(x, y, y', \dots, y^{(n)}) \\ f_2(x, y, y', \dots, y^{(n)}) \\ \dots \\ f_m(x, y, y', \dots, y^{(n)}) \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \dots \\ 0 \end{pmatrix}$$

- Example : SIR Model

- $\frac{dS}{dt} = -\beta IS, \frac{dI}{dt} = \beta IS - \nu I, \frac{dR}{dt} = \nu I$
 - $\mathbf{y} = [S, I, R], x = t$

- How to calculate \mathbf{y} for an arbitrary x



- Ordinary Differential Equation

- $$\begin{pmatrix} f_1(t, y, y', \dots, y^{(n)}) \\ f_2(t, y, y', \dots, y^{(n)}) \\ \vdots \\ f_m(t, y, y', \dots, y^{(n)}) \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

- How to calculate y for an arbitrary t

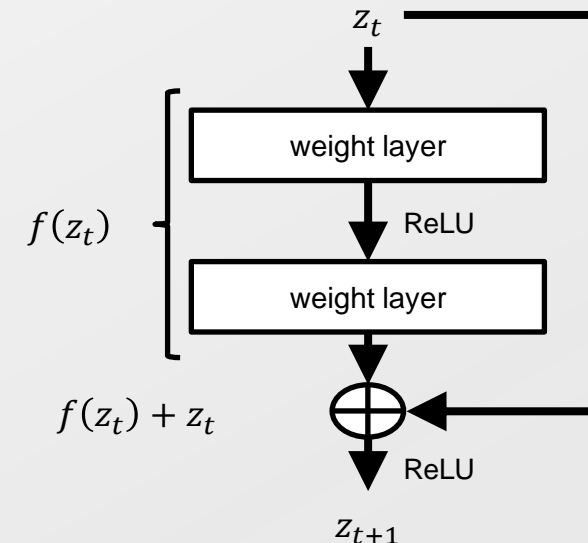
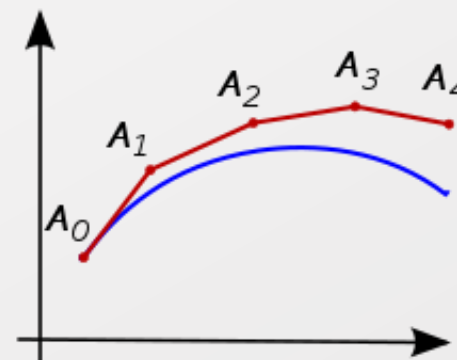
- Given an initial point and an ODE

- $y(t_0) = y_0$
- $y'(t) = f(t, y(t))$
- Let's assume a Taylor expansion on y at t_0 , and let's assume a very small step of h
 - $y(t_0 + h) = y(t_0) + hy'(t_0) + \frac{1}{2}h^2y''(t_0) + O(h^3)$
 - $y(t_0 + h) \approx y(t_0) + hf(t_0, y(t_0))$

- Euler method

- $y_{n+1} = y_n + hf(t_n, y_n)$
 - Assuming given $f(t_0, y_0)$
- can be expanded to Runge-Kutta method
 - Dynamic setting of h is feasible

https://en.wikipedia.org/wiki/Euler_method



$$F(z_t) = W_2 \sigma(W_1 z_t)$$

$$z_{t+1} = z_t + f(z_t; \{W_i\})$$

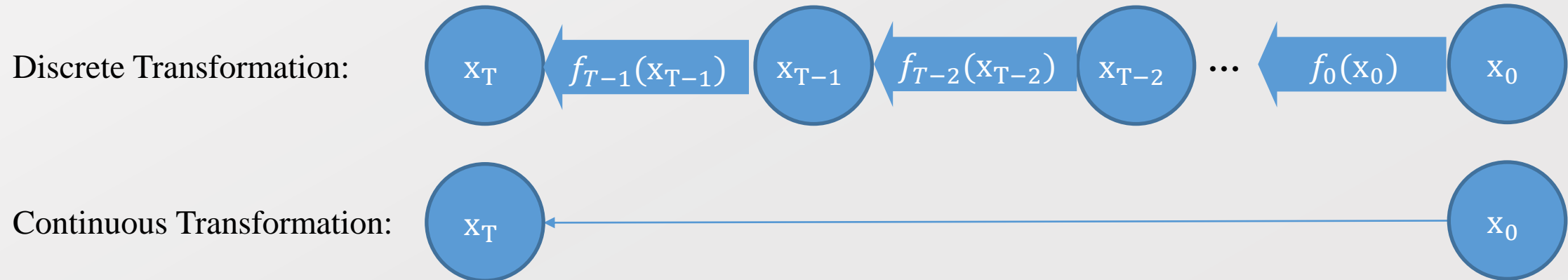
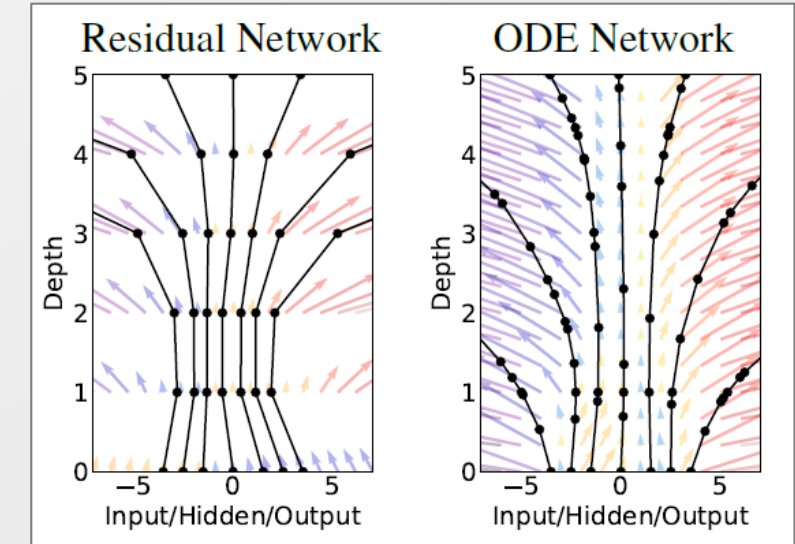
$$\rightarrow$$

$$z_{t+1} = z_t + \Delta t \frac{f(z_t; \theta)}{\Delta t} = z_t + \Delta t g(z_t, t; \theta)$$

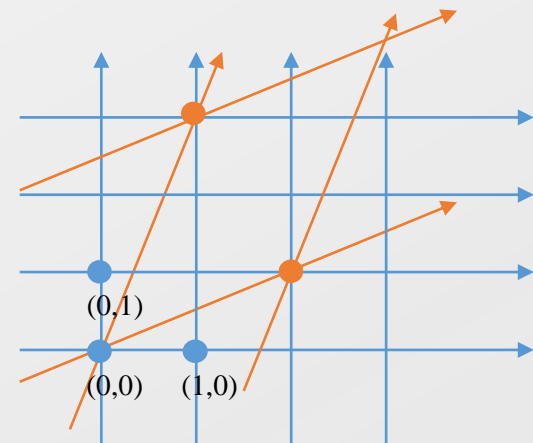
Exact Value Output vs. Differential Value Output from NN

- When we consider t as the layer number
 - Initially, t becomes the integer index of layer
 - However, is it feasible to see t as a continuous value?
- Usual neural network
 - $z_{t+1} = f_{NN_t}(z_t; \theta)$
 - z_{t+1} is the actual output given z_t
- Differential value integration with ODE solver
 - $\frac{dz_t}{dt} = g_{NN}(z(t), t; \theta)$
 - $z(t_1) = ODESolver(z(t_0), g_{NN}, t_0, t_1, \theta) = z(t_0) + (t_1 - t_0)g_{NN}(z(t_0), t_0; \theta)$
- Now, we can calculate z_t at t with precision values
- Then, how to view the flow model with the discrete transformation?

R. Chen et al., *Neural Ordinary Differential Equations*, NeurIPS 2018



- Compose the multiple transformation as a chain
 - $f = f_K \circ f_{K-1} \circ \dots \circ f_1$
 - Training with log-likelihood
 - $\max_{\theta} \left\{ \log p_Z(f_{\theta}(x)) + \sum_{i=1}^K \log \left| \det \frac{\partial f_i}{\partial f_{i-1}} \right| \right\}$
 - Role of $\det \frac{\partial f_i}{\partial f_{i-1}}$: Estimating the change of probability after the transformation
- Instantaneous Change of Variables
 - Assumption
 - Let $z(t)$ be a finite continuous random variable with $p(z(t))$ dependent on time
 - Let $\frac{dz}{dt} = \mathbf{g}(z(t), t)$ be a differential equation describing a continuous-in-time transformation of $z(t)$
 - Assuming that f is uniformly Lipschitz continuous in z and continuous in t
 - Then,
 - The change in log probability also follows a differential equation, $\frac{\partial \log p(z(t))}{\partial t} = -\text{tr} \left(\frac{d\mathbf{g}}{dz(t)} \right)$
 - \mathbf{g} does not have to be bijective, but \mathbf{g} should be Lipschitz continuous
- Continuous transformation in probability
 - $z(t_1) - z(t_0) = \int_{t_0}^{t_1} \mathbf{g}(z(t), t) dt$, where $z(t_1) = x$
 - $\log p(z(t_1)) - \log p(z(t_0)) = \int_{t_0}^{t_1} -\text{tr} \left(\frac{d\mathbf{g}}{dz(t)} \right) dt$, where $\log p(z(t_1)) = \log p(x)$
 - Log Determininant of Jacobian matrix \rightarrow Trace of Jacobian matrix
 - Bijective constraint \rightarrow Lipschitz constraint



- (Discrete) Normalizing Flow

- $\log p_{\theta}(x) = \log p(z_0) - \sum_{k=1}^K \log \det \left| \frac{\partial f_k}{\partial z_k} \right|$
- Complexity
 - $O(D^3)$ to calculate the log determinant

- Continuous Normalizing Flow

- $\log p_{\theta}(x) = \log p(z_0) - \int_{t_0}^{t_1} \text{tr} \left(\frac{dg}{dz(t)} \right) dt$
- Complexity
 - $O(D^2)$ to calculate the trace

- FFJORD: Continuous-time Flow with Hutchinson's trace estimator

- $\log p_{\theta}(x) = \log p(z_0) - \mathbb{E}_{p(\epsilon)} \left[\int_{t_0}^{t_1} v^T \left(\frac{dg}{dz(t)} \right) v dt \right]$
- Complexity
 - $O(D)$ to calculate the trace estimator

$$tr(A) = \mathbb{E}_{p(v)}[v^T A v]$$

Where $\mathbb{E}[v] = 0, Cov(v) = I$

$$\begin{pmatrix} \frac{\partial f_{\theta}^1}{\partial x_1} & \dots & \frac{\partial f_{\theta}^1}{\partial x_D} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_{\theta}^D}{\partial x_1} & \dots & \frac{\partial f_{\theta}^D}{\partial x_D} \end{pmatrix}$$