

Normalizing Flows (NF)

goals	autoencoder	VAE	GAN	NF
small codes		hyperparameter	$\dim(z)$	$\dim(z) = \dim(x)$
good reconstruction $\hat{X} = D(E(x)) \approx X$	✓	↑ (v) trade-off ↓ (v)	✗	✓ (some: perfect)
accurate distribution $p(x) \approx p^*(x)$	✗		✓	✓

idea of VF: generalize inverse transformation method to arbitrary dimensions

- $x \in \mathbb{R}^D$, $p^*(x)$ high-dim density, $A \subseteq \mathbb{R}^D$ a region

- define $\Pr(x \in A) = \int_A p^*(x) dx$

- let $z = f(x)$ an invertible encoding, $x = f^{-1}(z) := g(z)$

$\tilde{A} = f(A)$ the image of A in z -space $\tilde{A} = \{z : z = f(x) \text{ for } x \in A\}$

- consistency: for all A : $\Pr(z \in \tilde{A}) = \int_{\tilde{A}} q(z) dz \stackrel{!}{=} \Pr(x \in A)$ (*)

- apply the multi-dimensional change-of-variables formula of Analysis II

$$\int_{\tilde{A}=f(A)} q(z) dz = \int_{g(\tilde{A})} q(z=f(x)) |\det J_f| dx$$

Jacobian: $J_f = \begin{pmatrix} \frac{\partial f_1(x)}{\partial x_1} & \frac{\partial f_1(x)}{\partial x_2} & \dots & \frac{\partial f_1(x)}{\partial x_D} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_D(x)}{\partial x_1} & \dots & \dots & \frac{\partial f_D(x)}{\partial x_D} \end{pmatrix}$

- since consistency must hold for any A , the integrands must be equal
 \Rightarrow multi-variate change-of-variables formula

$$p(x) = q(z=f(x)) \left| \det J_f(x) \right|$$

- goal: - pre-define $q(z)$ (e.g. $q(z) = \mathcal{N}(0, \Pi)$)
 - learn $f(x)$, such that $p(x) \approx p^x(x)$ ($f(x) \Leftarrow$ neural network encoded)
- recall 1-D: $q(z) = \text{unif}(0,1) \Rightarrow f(x) = \text{CDF}_{p^x}(x)$
- for learning $q(z) = \mathcal{N}(0, \Pi)$ is better: - non-zero gradients for gradient descent
 - supported on all of \mathbb{R}^D
- two difficult problems in practice.

(1) how to ensure that $f(x)$ is invertible & efficiently compute $q(z) = f^{-1}(z)$

(2) how to (efficiently) calculate $\det J_f$

- start with problem (2)

2-D determinant easy. $\det \begin{pmatrix} a & b \\ c & d \end{pmatrix} = ad - bc$, 3D ok: $\det \begin{pmatrix} a & b & c \\ d & e & f \\ g & h & i \end{pmatrix} = a(ei - fh) + b(fg - di) + c(dh - eg)$

number of terms grows exponentially with D

- general solution with SVD

$$J = U \cdot \Lambda \cdot V^T \Rightarrow |\det J| = \det \Lambda = \prod_j \lambda_j$$

effort: $O(D^3)$

- exploit special case: if J is triangular

$$J = \begin{pmatrix} & & 0 \\ & \neq 0 & \\ & & \neq 0 \end{pmatrix}$$

$$\det J = \prod_j J_{jj}$$

we already talked about an instance: auto-regressive models

$$z_j = f_j(x_j, x_{<j})$$

$$J_f = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & 0 & 0 & \dots & 0 \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & 0 & \dots & 0 \\ \frac{\partial f_D}{\partial x_1} & \frac{\partial f_D}{\partial x_2} & \dots & \dots & \frac{\partial f_D}{\partial x_D} \end{pmatrix}$$

"triangular map"

"Kusche-Rosenblatt re-arrangement"

• if $f(x)$ is a multi-layer network, $f(x)$ is a composition of functions $f^{(l)}$

$$f = f^{(L)} \circ f^{(L-1)} \circ \dots \circ f^{(1)}$$

Jacobians of each layer

Jacobian of composition

$$J^{(l)} = \frac{\partial f^{(l)}(z^{(l-1)})}{\partial z^{(l-1)}}$$

$$J = \frac{\partial f(x)}{\partial x} = J^{(L)} \cdot J^{(L-1)} \cdot \dots \cdot J^{(1)}$$

$$z^{(l)} = f^{(l)}(z^{(l-1)}), \quad z^{(0)} = x, \quad z := z^{(L)}$$

$$\det J = \det J^{(L)} \cdot \dots \cdot \det J^{(1)}$$

\Rightarrow determinant of multi-layer network is easy when layer determinants are easy

\Rightarrow popular architecture: define all $f^{(l)}(z^{(l-1)})$ as auto-regressive functions

• relevant to problem (1): choose $f^{(l)}$ auto-regressive and easy to invert

major architecture: coupling layer \Rightarrow network "Real NVP" [Dinh et al 2017]

- idea: in each layer, change only half of the dimensions of $z^{(l-1)}$
pass the remaining dimensions or unchanged ("Skip connection")

$$z_j^{(e)} = z_j^{(e-1)} \quad \text{for } j=1, \dots, \tilde{D} \quad \tilde{D} = \lfloor D/2 \rfloor$$

$$z_j^{(e)} = f_j^{(e)} \left(z_j^{(e-1)}, \underbrace{z_{1:\tilde{D}}^{(e-1)}}_{\in \mathcal{Z}_{\tilde{D}}} \right) \quad \text{for } j = \tilde{D}+1, \dots, D$$

$$J^{(e)} = \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \\ \frac{\partial f_j^{(e)}}{\partial z_j^{(e-1)}} & \dots & \frac{\partial f_j^{(e)}}{\partial z_{\tilde{D}+1}^{(e-1)}} & 0 \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\partial f_j^{(e)}}{\partial z_j^{(e-1)}} & \dots & \frac{\partial f_j^{(e)}}{\partial z_{\tilde{D}}^{(e-1)}} & \frac{\partial f_j^{(e)}}{\partial z_j^{(e-1)}} \end{pmatrix} = \begin{pmatrix} I_{\tilde{D}} & 0 \\ \vdots & \vdots \\ \neq 0 & \text{diag} \left(\frac{\partial f_j^{(e)}}{\partial z_j^{(e-1)}} \right) \end{pmatrix}$$

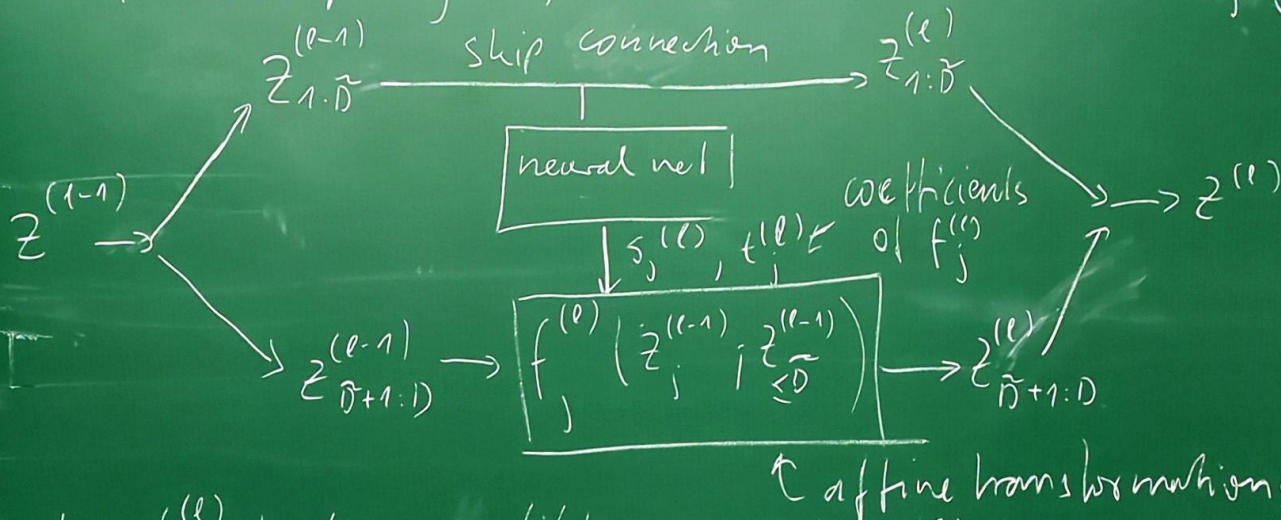
is triangular

$$\det J^{(e)} = \prod_{j=\tilde{D}+1}^D \frac{\partial f_j^{(e)}(z_j^{(e-1)}, z_{\leq \tilde{D}}^{(e-1)})}{\partial z_j^{(e-1)}}$$

• simplest invertible function: affine

$$z_j^{(l)} = s_j \cdot z_j^{(l-1)} + t_j$$

"affine coupling layer"



$$s_j^{(l)} \left(z_{\leq \tilde{D}}^{(l-1)} \right) \quad t_j^{(l)} \left(z_{\leq \tilde{D}}^{(l-1)} \right)$$

two neural nets →
(D, L in total)

• for $f_j^{(l)}$ to be invertible, we need $s_j^{(l)} \neq 0$

⇒ we actually learn $\tilde{s}_j^{(l)} \left(z_{\leq \tilde{D}}^{(l-1)} \right)$ and set $s_j^{(l)} = \exp(\tilde{s}_j^{(l)}) > 0$

(in practice, it is numerically more stable to set $s_j^{(l)} = \exp(\tanh(\tilde{s}_j^{(l)}))$)

• forward:

$$z_j^{(l)} = \exp\left(s_j^{(l)}(z_{\leq \tilde{D}}^{(l-1)})\right) \cdot z_j^{(l-1)} + t_j^{(l)}(z_{\leq \tilde{D}}^{(l-1)}) \quad j > \tilde{D}$$

$$z_j^{(l)} = z_j^{(l-1)} \quad j \leq \tilde{D}$$

• inverse:

$$z_j^{(l-1)} = z_j^{(l)} \quad \text{for } j \leq \tilde{D} \quad \text{thanks to skip w/ connection}$$

$$z_j^{(l-1)} = \left(z_j^{(l)} - t_j^{(l)}(z_{\leq \tilde{D}}^{(l-1)})\right) \cdot \exp\left(-s_j^{(l)}(z_{\leq \tilde{D}}^{(l-1)})\right) \quad \text{for } j > \tilde{D}$$

⇒ invertible layer constructed from non-invertible s and t (always executed in same direction)
 ⇒ choose s and t according to application (fully connected, convolutional, ...)

determinant of a affine coupling layer

$$\det J^{(l)} = \prod_{j=\tilde{D}+1}^D \frac{\partial f_j^{(l)}(z_j^{(l-1)}, z_{\leq \tilde{D}}^{(l-1)})}{\partial z_j^{(l-1)}} = \prod_{j=\tilde{D}+1}^D \exp\left(s_j^{(l)}(z_{\leq \tilde{D}}^{(l-1)})\right)$$