

Normalizing Flows (NFs)

- direct generalization of inverse transform sampling from 1-D to high dimensions

goal: $x \sim p^*(x)$, choose code distribution $q(z)$ ($= \mathcal{N}(0, I)$)

learn $g(z)$ such that $\hat{p}(x) = g_{\#} q(z) \approx p^*(x)$

$$\Leftrightarrow z \sim q(z) , \quad x = g(z) \sim \hat{p}(x)$$

[recall 1-D: $q(z) = \text{uniform}(0, 1)$ $f(x) = g^{-1}(x) = \text{CDF}(p^*(x))$]

- work-horse method of my team

generalize the change-of-variables formula to arbitrary dimension:

- 1-D: $p(x) = q(f(x)) |f'(x)|$

- D-dimensional: $p(x)$ with $\text{dom}(x)$ domain ($p(x) > 0$ if $x \in \text{dom}(x)$)

$A \subseteq \text{dom}(x)$ some subset, $\Pr(x \in A) = \int_A p(x) dx$

$z = f(x)$ and $p(z) = f_{\#} p(x)$, $\tilde{A} = \{z = f(x) : x \in A\} = f(A)$

transformation must preserve probability mass: $\Pr(z \in \tilde{A}) = \int_{\tilde{A}} p(z) dz$

$\int_{\tilde{A}=f(A)} p(z) dz = \int_{f^{-1}(\tilde{A})=A} p(z=f(x)) |\det f'(x)| dx$

$f'(x) = \begin{bmatrix} \frac{\partial z_1}{\partial x_1} & \frac{\partial z_2}{\partial x_1} & \dots & \frac{\partial z_D}{\partial x_1} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial z_1}{\partial x_D} & \frac{\partial z_2}{\partial x_D} & \dots & \frac{\partial z_D}{\partial x_D} \end{bmatrix}$ Jacobian matrix

to hold for all A.

$p(x) = p(z=f(x)) |\det f'(x)|$

general law about jacobians of invertible functions: $f(x)$, $g(z) = f^{-1}(z)$

$$f'(x) = (g'(z=f(x)))^{-1}, \text{ linear algebra: } \det B^{-1} = \frac{1}{\det B}$$

$$|\det f'(x)| = \frac{1}{|\det g'(z=f(x))|}$$

two equivalent forms of change-of-variables formula

$$\begin{aligned} p(X=x) &= q(z=f(x)) \cdot |\det f'(x)| \\ p(X=g(z)) &= q(z=z) \cdot |\det g'(z)|^{-1} \end{aligned}$$

\Rightarrow NF learning problem: ① How to define $f(x)$ such that $g(z)=f^{-1}(z)$ exists and is easy to calculate?
② How can we efficiently calculate $|\det f'(x)|$ and/or $|\det g'(z)|$?
③ How to learn $f(x)$ and/or $g(z)$?

calculating a determinant: $D=1$: $\det f'(\lambda) = f'(x)$ because $f'(x)$ is scalar

$$D=2 \quad \det \begin{pmatrix} a & b \\ c & d \end{pmatrix} = a \cdot d - b \cdot c$$

$$D=3 \quad \det \begin{pmatrix} a & b & c \\ d & e & f \\ g & h & i \end{pmatrix} = a \cdot (e \cdot i - f \cdot h) + b \cdot (f \cdot g - d \cdot i) + c \cdot (d \cdot h - e \cdot g)$$

number of terms grows fast with $D \Rightarrow$ impractical

solution via SVD: $f'(x) = U \cdot \Lambda \cdot V^T$

U : orthogonal
 Λ : diagonal, ≥ 0

still too expensive for training (millions of evaluations)

solution in NF based on triangular maps: $\det \left[B = \begin{pmatrix} b_{11} & 0 & 0 \\ b_{21} & b_{22} & 0 \\ \vdots & \vdots & \vdots \\ b_{D1} & 0 & b_{DD} \end{pmatrix} \right] = \prod_{j=1}^D b_{jj}$

$$|\det f'(x)| = |\det U| \cdot |\det \Lambda| \cdot |\det V^T|$$
$$= 1 \cdot \prod_{j=1}^D \lambda_{jj} \cdot 1$$
$$= \prod_{j=1}^D \lambda_{jj}$$

to make $f'(x)$ triangular, $f(x)$ must be a triangular map ("Leath-Rosenblatt-rearrangement")

$$f(x) = \begin{bmatrix} f_1(x) \\ f_2(x) \\ \vdots \\ f_D(x) \end{bmatrix} = \begin{bmatrix} f_1(x_1) \\ f_2(x_1, x_2) \\ \vdots \\ f_D(x_1, \dots, x_D) \end{bmatrix} \Rightarrow \frac{\partial^2 j}{\partial x_{j'}} = \frac{\partial f_j(x)}{\partial x_{j'}} = 0 \quad \text{if } j' > j \Rightarrow \text{upper triangle of } f'(x) \text{ is zero}$$

(a abbreviation: $f_j(x_{\leq j})$)

$$\Rightarrow \text{def } f'(x) = \left\| \frac{\partial f_j(x)}{\partial x_j} \right\|_{j=1}^D \leftarrow \text{diagonal elements of } f'(x)$$

However, learning a triangular map is hard & impractical \Rightarrow NF use a compromise.

$f(x) = f_L \circ f_{L-1} \circ \dots \circ f_2 \circ f_1$ is a composition of L layers, and the f_ℓ are triangular maps

$$|\det f'(x)| = |\det f'_L| \circ |\det f'_{L-1}| \dots |\det f'_1|$$

choose f^l such that - $\det f'_e$ is cheap

all easy because triangular maps

- f_e^{-1} is cheap

$$\rightarrow g(z) = \underbrace{g_1 \circ}_{f_1^{-1}} \underbrace{g_2 \circ}_{f_2^{-1}} \dots \underbrace{g_L(z)}_{f_0^{-1}}$$

- $f(x)$ is expressive

(can learn complicated $p^*(x)$)

coupling layer
("affine" coupling layer)

$$z^{(e)} = f_e(z^{(e-1)}) \quad \text{with } z^{(0)} = x$$

split vector $z^{(e)}$ into two halves

$$z^{(e)} = \begin{bmatrix} z_{\leq \tilde{D}}^{(e)} \\ z_{> \tilde{D}}^{(e)} \end{bmatrix} \quad \text{intermediate activations} \quad \tilde{D} = \lfloor \frac{D}{2} \rfloor$$

$$(1) \quad z_{\leq \tilde{D}}^{(e)} = z_{\leq \tilde{D}}^{(e-1)} \quad (\text{skip connection for first half})$$

$$(2) \quad z_i^{(e)} = s_i^{(e)}(z_{\leq \tilde{D}}^{(e-1)}) \cdot z_i^{(e-1)} + t_i^{(e)}(z_{\leq \tilde{D}}^{(e-1)}) \quad i > \tilde{D} \quad (\text{second half})$$

$$z_{\leq \tilde{D}}^{(e)} = \begin{bmatrix} z_1^{(e)} \\ \vdots \\ z_{\tilde{D}}^{(e)} \end{bmatrix} \quad z_{> \tilde{D}}^{(e)} = \begin{bmatrix} z_{\tilde{D}+1}^{(e)} \\ \vdots \\ z_D^{(e)} \end{bmatrix}$$

$z^{(e)} = f_e(z^{(e-1)})$ defined a) a affine coupling is a triangular map

if $j \leq \tilde{n}$: $\frac{\partial z_j^{(e)}}{\partial z_{j'}^{(e-1)}} = \begin{cases} 1 & \text{if } j=j' \\ 0 & \text{otherwise} \end{cases}$ (skip connections)

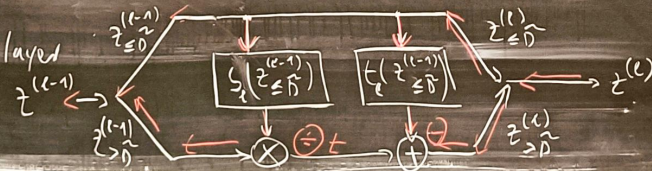
if $j > \tilde{n}$: $\frac{\partial z_j^{(e)}}{\partial z_{j'}^{(e-1)}} = \begin{cases} 0 & \text{if } j' > j \\ s(z_{\leq \tilde{n}}^{(e-1)}) & \text{if } j=j' \\ \neq 0 & \text{otherwise} \end{cases}$

$|\det f'(z^{(e-1)})| = \prod_{j=\tilde{n}+1}^D s_j(z_{\leq \tilde{n}}^{(e-1)})$

forward: $z_{>\tilde{n}}^{(e)} = s_e(z_{\leq \tilde{n}}^{(e-1)}) \cdot z_{>\tilde{n}}^{(e-1)} + t_e(z_{\leq \tilde{n}}^{(e-1)})$

$f'(z^{(e-1)}) = \begin{bmatrix} 1 & & 0 & & \\ & \ddots & & & \\ & & 1 & & \\ & & & \ddots & \\ & & & & 1 & & \\ & & & & & s_{\tilde{n}+1}(z_{\leq \tilde{n}}^{(e-1)}) & 0 \\ & & & & & & \ddots & \\ & & & & & & & 0 & s_D(z_{\leq \tilde{n}}^{(e-1)}) \end{bmatrix}$

picture of layer



require: $s_e(z_{\leq \tilde{n}}^{(e-1)}) \neq 0$

inverse: $z_{>\tilde{n}}^{(e-1)} = \frac{z_{>\tilde{n}}^{(e)} - t_e(z_{\leq \tilde{n}}^{(e)})}{s_e(z_{\leq \tilde{n}}^{(e)})}$