

Deep Learning-based approach to

(1) Density est. (2) sampling from dists

Setting: $X_1, \dots, X_n \sim P_X$

We would like to (1) estimate the density P_X

(2) sample from P_X

Approaches:

deep learning
based

(1) Parametric (mixture models)

(2) Energy-based models (markov random fields)

(3) VAE (variational autoencoders) Kingma, Welling 2013

(4) GAN (generative adversarial networks) Goodfellow et al. 2014

(5) Normalizing Flows

(6) Diffusion

	(1) MC	(2) uncond density	(3) M	(4) H	(5) ✓	(6) J
(1) ability to sample	✓	✓		✓	✓	
(2) evaluate density	✓	X	X	X	✓	J
(3) expressiveness dimensions Scalability	M	X	M	H		✓
(4) training/ estimate/ learn	MC	ML	SGD ELBO	GEP game theory		SGD

Normalizing flows

Approach:

(1) Select a r.v. Z that has a density that is easy to evaluate and whose distn is easy to sample from.

(2) Assume there exists an invertible (bijection) differentiable function f_θ , s.t. $Z = f_\theta(X)$ $f: \mathcal{X} \rightarrow \mathcal{Z}$

(3) Find $\hat{\theta} = \operatorname{argmax}_{\theta \in \Theta} \log(L(D, \theta))$ (where $X \sim P_X$, $Z \sim P_Z$)

$$\begin{aligned} &= \operatorname{argmax}_{\theta \in \Theta} \log \left(\prod_{i=1}^n p_Z(x_i) \right) \\ &= \operatorname{argmax}_{\theta \in \Theta} \sum_{i=1}^n \log(p_Z(x_i)) \end{aligned}$$

★

discussed
on page 13



now is this a func of θ ?

Recall: Change of variable formula

Given $\mathcal{X} \sim P_X$, $\mathcal{Z} \sim P_Z$, and $f: \mathcal{X} \xrightarrow{P} \mathcal{Z}$

where f is a bijection \Rightarrow differentiable, then

$$p_Z(x) = p_X(f(x)) |\det(J_f(x))|$$

where $J_f(x)$ is the Jacobian of f evaluated at x
e.g. Matrix of partial derivatives $f: \mathbb{R}^p \rightarrow \mathbb{R}^q$

Note: $\det(X)$ is only defined for square matrices so $X \in \mathbb{R}^{n \times n}$
must have the same dim.

$$\begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_p} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_q}{\partial x_1} & \cdots & \frac{\partial f_q}{\partial x_p} \end{bmatrix}$$

(cont.)

$$\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmax}} \sum_{i=1}^n \log(p_{\mathcal{X}}(x_i))$$

$$\stackrel{\text{change of var}}{=} \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^n \log \left[p_{\mathcal{Z}}(f_{\theta}(x_i)) |\det(J_{f_{\theta}}(x_i))| \right]$$

$$= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^n \log \left[p_{\mathcal{Z}}(f_{\theta}(x_i)) \right] + \sum_{i=1}^n \log(\det J_{f_{\theta}}(x_i))$$

↑ fit with gradient descent

$$\sum_{k=1}^n \log[\det(J_{f_{\theta_k}}(x_i))]$$

(4) (A) To evaluate the density of a new point:

$$p_{\mathcal{X}}(x_j) = p_{\mathcal{Z}}(f_{\theta}(x_j)) |\det(J_{f_{\theta}}(x_j))|$$

(B) To sample from $p_{\mathcal{X}}$,

(1) sample from $p_{\mathcal{Z}}$ to obtain z_j

(2) $x_j = f_{\theta}^{-1}(z_j)$ & why invertibility is (practically) needed

Terminology: f is called a "flow" (as it flows $P_Z \rightarrow P_{Z'}$) and since Z' is usually chosen to be $\sim N$, \rightarrow normalizing flows

Takeaways:

- unusual application of C.o.V. ("backwards")
- don't directly estimate the quantity / function of interest. Instead learn a map from it to a quantity / function you know.

Potential Trouble

(1) dimensionality of f (must be $Z^{P_Z} \rightarrow Z^{P_{Z'}}$) may cause issues for:

(1A) f^{-1}_θ and (1B) $\det(J_{f_\theta}(x_i))$

(2) need f to be very (arbitrarily) "flexible"

Next: how to constraint such f 's

Starting with "programmatic" approach to achieving (2)

If would be nice to break f into simple building blocks. (so that we can build arbitrarily complex flows from simple building blocks)

Recall: invertible, differentiable functions are closed under composition.

If f_j are invertible & diff. then

$$f(x) = f_k \circ f_{k-1} \circ \dots \circ f_1(x) \quad (\star)$$

is also, in particular

$$f^{-1}(z) = f_1^{-1} \circ f_2^{-1} \dots \circ f_k^{-1}(z)$$

Recall: $\det(J_f) = \det\left(\prod_{k=1}^K J_{f_k}\right) = \prod_{k=1}^K \det(J_{f_k})$
(for f of the type defined in \star)

→ so it is enough to have an easy-to-calc $\det(J_{f_k})$

Now, we consider how to construct f that address (1A); (1B) starting with simplest option:

Linear flow

$$z = f_k(x) = A_k x + b_k$$

$$\text{inverse: } f_k^{-1}(z) = A_k^{-1}(z - b_k)$$

$$\det(J_f) = \det(A_k)$$

issues: (1) linear fns are closed under composition.

(2) inv and det can be computationally hard depending on structure of A_k

What kinds of f have easy f_k^{-1} and $\det(J_f)$?

<u>A_k</u>	inverse (f_k^{-1}) (A_k)	$\det(J_f)$ $(\det(A))$
full	$O(p^3)$	$O(p^3)^*$
diag	$O(p)$	$O(p)$
triangular	$O(p^2)$	$O(p)$
block diag	$O(C^3 p)$	$O(C^3 p)$

Recall:

(A) The determinant of a triangular matrix is the product of the diagonals (can see via Laplace expansion)

$$J = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & & & 0 \\ & \ddots & & \\ & & \ddots & 0 \\ \frac{\partial f_p}{\partial x_1} & & & \frac{\partial f_p}{\partial x_p} \\ & \ddots & & \ddots & \ddots \end{bmatrix}$$

$$\det(J) = \prod_{j=1}^p \frac{\partial f_j}{\partial x_j}$$

(B) For a block triangular matrix,

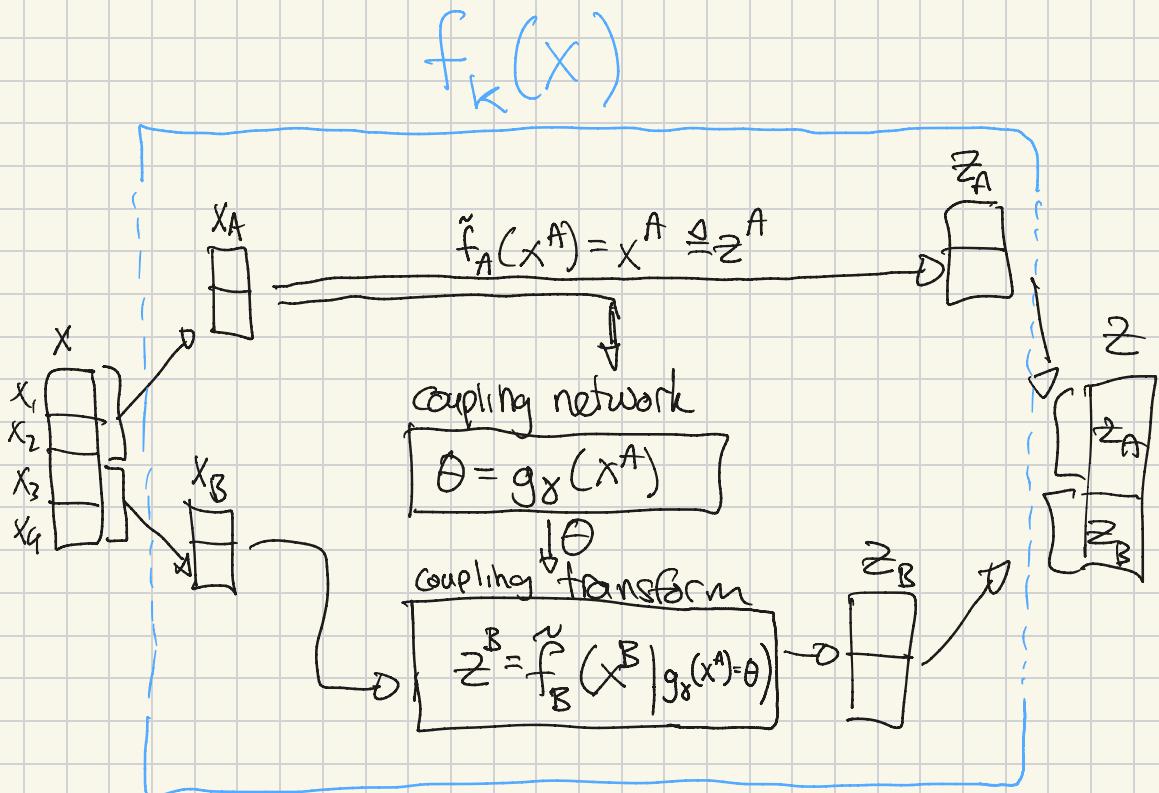
$$J = \begin{bmatrix} A & 0 \\ C & B \end{bmatrix} \quad \det(J) = \det(A) \det(B)$$

Now, how to build a flow (f_k) that is more expressive than the linear flow while also having a "simple" Jacobian?

Next, we look at an example of such a flow (whose Jacobian has a block triangular structure)

Coupling flow

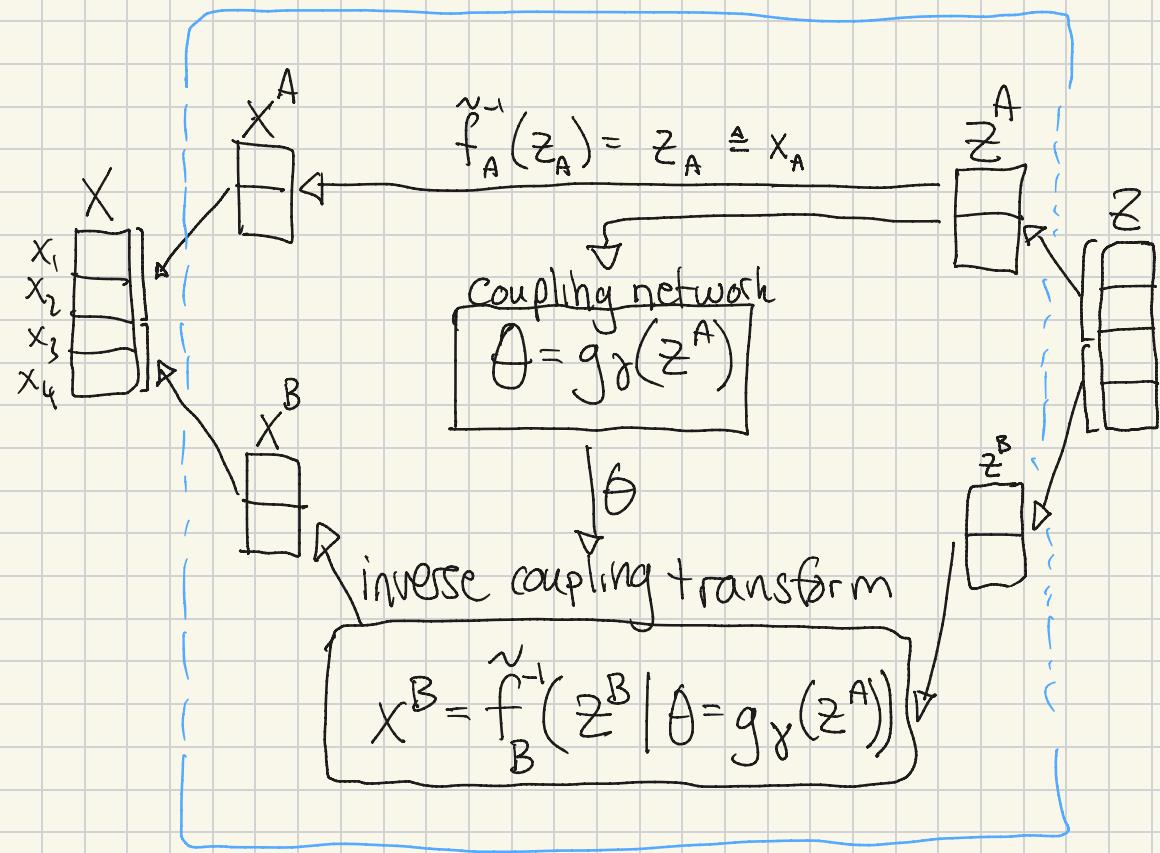
- This flow (f_k) has a general structure that results in a "simple" $\det(J_{f_k})$ while also allowing f_k to capture complex dependencies (via g_θ introduced below)
- Next we will define f_k , f_k^{-1} , and J_{f_k}
- First, we define the "forward" flow ($f_k(x)$)



$$f_k(x) = \begin{bmatrix} x^A = \tilde{f}_A(x_A) \\ \tilde{f}_B(x^B | \Theta = g_\theta(x^A)) \end{bmatrix} \triangleq \begin{bmatrix} f_k^A(x) \\ f_k^B(x) \end{bmatrix}$$

- Next we define the inverse $f_k^{-1}(z)$:

$$f_k^{-1}(z)$$



$$f_k^{-1}(z) = \begin{bmatrix} z_A = \tilde{f}_A^{-1}(z_A) \\ \tilde{f}_B^{-1}(z^B | \theta = g_\delta(z^A)) \end{bmatrix}$$

o Finally, we define the Jacobian (J_{f_k}):

$$J_{f_k}(x) = \begin{bmatrix} \frac{\partial f_k^A}{\partial x_A} & \frac{\partial f_k^A}{\partial x_B} \\ \frac{\partial f_k^B}{\partial x_A} & \frac{\partial f_k^B}{\partial x_B} \end{bmatrix} = \begin{bmatrix} I & II \\ III & IV \end{bmatrix}$$

$$= \begin{bmatrix} I \\ \frac{\partial}{\partial x_A} \left[\tilde{f}_B(x^B | g(x^A) = \theta) \right] \end{bmatrix} \quad \begin{bmatrix} \textcircled{I} \\ \frac{\partial}{\partial x_B} \left[\tilde{f}_B(x^B | g(x^A) = \theta) \right] \end{bmatrix}$$

In the
IV partial
derivative,
 θ is constant
(wrt x_B)

$$= \begin{bmatrix} I \\ \frac{\partial}{\partial x_A} \left[\tilde{f}_B(x^B | g(x^A) = \theta) \right] \end{bmatrix} \quad \begin{bmatrix} \textcircled{I} \\ \frac{\partial}{\partial x_B} \left[\tilde{f}_B(x^B | \theta) \right] \end{bmatrix}$$

$$= \begin{bmatrix} I \\ \frac{\partial}{\partial x_A} \left[\tilde{f}_B(x^B | g(x^A) = \theta) \right] \end{bmatrix} \quad \begin{bmatrix} \textcircled{I} \\ J_{\tilde{f}_B}(x^B | \theta) \end{bmatrix}$$

III might be complex but is ignored in the determinant's

Some important properties of f_k :

(1) Invertibility: For f_k^{-1} to be simple, it is sufficient for \tilde{f}_B to be easy to invert (w.r.t $\underline{x_B}$)

(2) The determinant of the Jacobian (J_{f_n}):

J_{f_n} is block (lower) triangular, so

$$\det(J_{f_n}) = \det(I) \cdot \det(J_{\tilde{f}_B}(x^B | \theta)) = \det(J_{\tilde{f}_B}(x^B | \theta))$$

Thus, for f_n to have a tractable $\det(J_{f_n})$, it is sufficient for $\det(J_{\tilde{f}_B}(x^B | \theta))$ to be tractable.

This will be the case if \tilde{f}_B is a simple function wrt x_B (not necessarily wrt x_A).

(3) Expressivity of f_k : Notice that no

restrictions are needed on the complexity/form of $g_\theta(x^A)$ (to produce easy f_k^{-1} and $\det(J_{f_n})$).

Thus, \tilde{f}_B can be made arbitrarily complex (ie nonlinear) as a function of x_A through its dependence on $\theta = g_\theta(x_A)$ via the choice of this g_θ (e.g. a DNN or similar)

What are some choices for \tilde{f}_B that are easy to invert and have simple $\det(J_{\tilde{f}_B})$?

- Additive (NICE, Dinh et al. 2014)

$$\tilde{f}_B(x_B | t(x_A)) = x_B + t(x_A)$$

$$\tilde{f}_B^{-1}(z_B | t(z_A)) = z_B - t(z_A)$$

$$\det(J_{\tilde{f}_B}) = I$$

- Affine (RealNVP, Dinh et al. 2016)

$$\tilde{f}_B(x_B | t(x_A), s(x_A)) = \exp\{s(x_A)\} \odot x_B + t(x_A)$$

$$\tilde{f}_B^{-1}(z_B | t(z_A), s(z_A)) = \frac{(z_B - t(z_A))}{\exp\{s(z_A)\}}$$

$$\det(J_{\tilde{f}_B}) = \det\left(\begin{bmatrix} \exp\{s_1(x_A)\} & & & \\ & \ddots & & \\ & & \textcircled{O} & \\ & & & \exp\{s_q(x_A)\} \end{bmatrix}\right)$$

$$= \exp\left\{\sum_{j=1}^q s_j(x_A)\right\}$$

$$q = \dim(x_B)$$

$= \frac{P}{2}$ for the f_k
we diagrammed

Now we have constructed an f_n that is expressive (via g_x), with tractable f_n^{-1} and $\det(J_{f_n})$ (via \tilde{f}_B).

One outstanding issue remains, namely that directly composing f_n of the type we defined would leave x_A unchanged i.e.

$$\begin{bmatrix} x_A \\ z_B \end{bmatrix} = f(x) = (f_K \circ \dots \circ f_1)(x)$$

One solution is to alternate f_n with random (\neq fixed) permutations (p_n), thus:

$$f(x) = (f_K \circ p_K \circ \dots \circ f_1 \circ p_1)(x)$$

Note that permutations are invertible & differentiable, meeting the requirements for the change of variable formula, and they have simple inverse and $\det(J) (\pm 1)$.

Summary

Developed a general structure for flows that is expressive with simple $f^{-1} \in \text{det}(J_f)$

One important outstanding issue. We began with the critical assumption that $\exists f_0 \in F$ s.t. $Z \stackrel{d}{=} f_0(X)$ where $X \sim P_X, Z \sim P_Z$ but then we restricted our search space for $f_0(F)$ to a very specific function family (coupling flows). So,

- (1) Is there any guarantee that f_0 exists such that $Z \stackrel{d}{=} f_0(X)$ for this particular (P_X, P_Z) pair?
- (2) If yes, then is $f_0 \in F$ also guaranteed? (i.e. can we be sure that one parametrization of the coupling flow is this f_0 ?)

Key result for (1) :

(1) Existence (Lemma 2.1, Bogachev et al 2005)

(roughly) For absolutely continuous

P_x, P_z on \mathbb{R}^P , $\exists!$ increasing, triangular mapping f_θ s.t. $Z \stackrel{d}{=} f_\theta(X)$. Where

- a mapping $f_\theta(x^{(1)}, \dots, x^{(P)}) \triangleq (f_\theta^{(1)}, \dots, f_\theta^{(P)})$ is

triangular if $\forall j$ $f_\theta^{(j)}$ is a function of

$(x^{(1)}, \dots, x^{(j)})$ (and thus has triangular

Jacobian)

- a mapping $f_\theta(x^{(1)}, \dots, x^{(P)}) \triangleq (f_\theta^{(1)}, \dots, f_\theta^{(P)})$

is increasing if $f_\theta^{(j)}$ is a (strictly) increasing function w.r.t $x^{(j)}$ $\forall j$

[see also Knöthe-Rosenblatt map/rearrangement]

Universality: Regarding (2) (ie whether $f \in \mathcal{F}$), this (unsurprisingly) depends on \mathcal{F} .
Some classes of \mathcal{F} for which this universality has been shown :

- neural autoregressive flows (Huang et al, 2018)
- augmented neural ODE (Dupont et al, 2019)
- sum of squares polynomial flow (Jain et al 2019)

[see also a more detailed survey of the]
Statistical properties of triangular flows : "Triangular Flows for Generative Modeling - Statistical Consistency, Smoothness Classes, and Fast Rates" Irons et al 2022