

[Paper review 48]

Normalizing Flows for Probabilistic Modeling and Inference

(Papamakarios, et al., 2019)

[Contents]

1. Abstract

NF : provide expressive distn, require 2 things

- (1) base distn
- (2) series of bijective transformation

Provide a perspective by describing flows through the lens of **probabilistic modeling and inference**

2. Introduction

How is the data generated (produced) ?

Build probability distn as NF!

- Section 2) formal & conceptual structure of NF
- Section 3) in detail for finite & infinitesimal variants
- Section 4) general perspective
- Section 5) extensions to structured domains & geometries
- Section 6) oft-encountered applications

3. Normalizing Flow

3-1. Definition & Basics

$\mathbf{x} = T(\mathbf{u})$ where $\mathbf{u} \sim p_{\mathbf{u}}(\mathbf{u})$

- $p_{\mathbf{u}}(\mathbf{u})$: base distn
- params
 - ϕ : for transformation T
 - θ : for base distn $p_{\mathbf{u}}(\mathbf{u})$

- T : must be **invertible** & **differentiable**

T^{-1} : must be **differentiable**

→ such T are called **diffeomorphisms**

Change of variables

- $p_{\mathbf{x}}(\mathbf{x}) = p_{\mathbf{u}}(\mathbf{u}) |\det J_T(\mathbf{u})|^{-1}$ where $\mathbf{u} = T^{-1}(\mathbf{x})$.
($= p_{\mathbf{x}}(\mathbf{x}) = p_{\mathbf{u}}(T^{-1}(\mathbf{x})) |\det J_{T^{-1}}(\mathbf{x})|$)

Jacobian $J_T(\mathbf{u})$

- is the $D \times D$ matrix of all partial derivatives of T

$$J_T(\mathbf{u}) = \begin{bmatrix} \frac{\partial T_1}{\partial u_1} & \cdots & \frac{\partial T_1}{\partial u_D} \\ \vdots & \ddots & \vdots \\ \frac{\partial T_D}{\partial u_1} & \cdots & \frac{\partial T_D}{\partial u_D} \end{bmatrix}.$$

Typically, $T : \text{NN} \& p_u(\mathbf{u}) : \text{MVN}$

Absolute Jacobian determinant : $|\det J_T(\mathbf{u})|$

- quantifies the relative change of volume of a small neighbourhood around \mathbf{u} due to T .

Property of "invertible" + "differentiable" transformation = **composable**

$$(T_2 \circ T_1)^{-1} = T_1^{-1} \circ T_2^{-1}$$

$$\det J_{T_2 \circ T_1}(\mathbf{u}) = \det J_{T_2}(T_1(\mathbf{u})) \cdot \det J_{T_1}(\mathbf{u}).$$

"Normalizing" "Flow"

- (1) Flow : trajectory, that a collection of samples from $p_{\mathbf{u}}(\mathbf{u})$ follow
by the sequence of transformations T_1, \dots, T_K
- (2) Normalizing : inverse flow through $T_K^{-1}, \dots, T_1^{-1}$ takes a collection of samples from $p_{\mathbf{x}}(\mathbf{x})$
transforms it back (=normalize them) into a collection of samples from $p_{\mathbf{u}}(\mathbf{u})$

Flow-based model provides 2 operations

- **(1) SAMPLING from the model** (forward transformation)
→ $\mathbf{x} = T(\mathbf{u})$ where $\mathbf{u} \sim p_{\mathbf{u}}(\mathbf{u})$
- **(2) EVALUATING the model's density** (Inverse transformation & Jacobian determinant)
→ $p_{\mathbf{x}}(\mathbf{x}) = p_{\mathbf{u}}(T^{-1}(\mathbf{x})) |\det J_{T^{-1}}(\mathbf{x})|$

3-2. Expressive power of flow-based models

How expressive is it?

→ **universal representation is possible, under reasonable conditions on $p_x(\mathbf{x})$**

3-3. Using flows for modeling & inference

how to fit **flow-based model** $p_{\mathbf{x}}(\mathbf{x}; \theta)$ to a **target distn** $p_{\mathbf{x}}^*(\mathbf{x})$?

- by minimizing some divergence/discrepancy (ex. KL-div)
- params : $\theta = \{\phi, \psi\}$
where ϕ are the parameters of T and ψ are the parameters of $p_{\mathbf{u}}(\mathbf{u})$

3-3-1. Forward KL & MLE

$$\begin{aligned}\mathcal{L}(\theta) &= D_{\text{KL}} [p_{\mathbf{x}}^*(\mathbf{x}) \| p_{\mathbf{x}}(\mathbf{x}; \theta)] \\ &= -\mathbb{E}_{p_{\mathbf{x}}^*(\mathbf{x})} [\log p_{\mathbf{x}}(\mathbf{x}; \theta)] + \text{const.} \\ &= -\mathbb{E}_{p_{\mathbf{x}}^*(\mathbf{x})} [\log p_{\mathbf{u}}(T^{-1}(\mathbf{x}; \phi); \psi) + \log |\det J_{T^{-1}}(\mathbf{x}; \phi)|] + \text{const.}\end{aligned}$$

- suitable when we have **samples from the target distn**
- unsuitable when we cannot evaluate target density $p_{\mathbf{x}}^*(\mathbf{x})$

Using MC samples...

$$\mathcal{L}(\theta) \approx -\frac{1}{N} \sum_{n=1}^N \log p_{\mathbf{u}}(T^{-1}(\mathbf{x}_n; \phi); \psi) + \log |\det J_{T^{-1}}(\mathbf{x}_n; \phi)| + \text{const.}$$

Minimizing KL-div

= Fitting flow-based model to the samples $\{\mathbf{x}_n\}_{n=1}^N$ **by MLE**

How to optimize? **Iteratively with SGD**

$$\begin{aligned}\nabla_{\phi} \mathcal{L}(\theta) &\approx -\frac{1}{N} \sum_{n=1}^N \nabla_{\phi} \log p_{\mathbf{u}}(T^{-1}(\mathbf{x}_n; \phi); \psi) + \nabla_{\phi} \log |\det J_{T^{-1}}(\mathbf{x}_n; \phi)| \\ \nabla_{\psi} \mathcal{L}(\theta) &\approx -\frac{1}{N} \sum_{n=1}^N \nabla_{\psi} \log p_{\mathbf{u}}(T^{-1}(\mathbf{x}_n; \phi); \psi)\end{aligned}$$

- $\nabla_{\psi} \mathcal{L}(\theta)$: can be done in closed-form if $p_{\mathbf{u}}(\mathbf{u}; \psi)$ admits closed-form MLE
(ex. Gaussian distn)

Fitting by **MLE**

- need to compute T^{-1} , Jacobian determinant, density $p_{\mathbf{u}}(\mathbf{u}; \psi)$
& different those three
- no need to compute T , or sample from $p_{\mathbf{u}}(\mathbf{u}; \psi)$

3-3-2. Reverse KL

$$\begin{aligned}\mathcal{L}(\theta) &= D_{\text{KL}} [p_{\mathbf{x}}(\mathbf{x}; \theta) \| p_{\mathbf{x}}^*(\mathbf{x})] \\ &= \mathbb{E}_{p_{\mathbf{x}}(\mathbf{x}; \theta)} [\log p_{\mathbf{x}}(\mathbf{x}; \theta) - \log p_{\mathbf{x}}^*(\mathbf{x})] \\ &= \mathbb{E}_{p_{\mathbf{u}}(\mathbf{u}; \psi)} [\log p_{\mathbf{u}}(\mathbf{u}; \psi) - \log |\det J_T(\mathbf{u}; \phi)| - \log p_{\mathbf{x}}^*(T(\mathbf{u}; \phi))]\end{aligned}$$

- suitable when we can evaluate target density $p_{\mathbf{x}}^*(\mathbf{x})$
- no need to **samples from the target distn**

Rewrite reverse KL

- $p_{\mathbf{x}}^*(\mathbf{x}) = \tilde{p}_{\mathbf{x}}(\mathbf{x}) / C$ (where $\tilde{p}_{\mathbf{x}}(\mathbf{x})$ is unnormalized density)
- $\mathcal{L}(\theta) = \mathbb{E}_{p_{\mathbf{u}}(\mathbf{u}; \psi)} [\log p_{\mathbf{u}}(\mathbf{u}; \psi) - \log |\det J_T(\mathbf{u}; \phi)| - \log \tilde{p}_{\mathbf{x}}(T(\mathbf{u}; \phi))] + \text{const.}$

Minimize reverse KL iteratively with SGD & Reparam Trick & MC estimation

- $\{\mathbf{u}_n\}_{n=1}^N$: samples from $p_{\mathbf{u}}(\mathbf{u}; \psi)$
- minimize w.r.t ϕ :
$$\nabla_{\phi} \mathcal{L}(\theta) \approx -\frac{1}{N} \sum_{n=1}^N \nabla_{\phi} \log |\det J_T(\mathbf{u}_n; \phi)| + \nabla_{\phi} \log \tilde{p}_{\mathbf{x}}(T(\mathbf{u}_n; \phi)).$$
- minimize w.r.t ψ :
use reparameterization ... $\mathbf{u} = T'(\mathbf{u}'; \psi)$ where $\mathbf{u}' \sim p_{\mathbf{u}'}(\mathbf{u}')$
but we can absorb T' into T & replace the base distn with $p_{\mathbf{u}'}(\mathbf{u}')$

3-3-3. Duality between Forward & Reverse KL

Think of

- target $p_{\mathbf{x}}^*(\mathbf{x})$ as "base distribution"
- inverse flow as "inducing a distn $p_{\mathbf{u}}^*(\mathbf{u}; \phi)$ "

$p_{\mathbf{u}}^*(\mathbf{u}; \phi) = p_{\mathbf{u}}(\mathbf{u}; \psi)$ if and only if $p_{\mathbf{x}}^*(\mathbf{x}) = p_{\mathbf{x}}(\mathbf{x}; \theta)$ Thus (a) & (b) are equivalent!

- (a) fitting model $p_{\mathbf{x}}(\mathbf{x}; \theta)$ to target $p_{\mathbf{x}}^*(\mathbf{x})$
- (b) fitting induced distn $p_{\mathbf{u}}^*(\mathbf{u}; \phi)$ to base $p_{\mathbf{u}}(\mathbf{u}; \psi)$

Using change of variables... $D_{\text{KL}} [p_{\mathbf{x}}^*(\mathbf{x}) \| p_{\mathbf{x}}(\mathbf{x}; \theta)] = D_{\text{KL}} [p_{\mathbf{u}}^*(\mathbf{u}; \phi) \| p_{\mathbf{u}}(\mathbf{u}; \psi)]$.

→ "fitting the model to the target using the **forward KL divergence**

= fitting the induced distribution $p_{\mathbf{u}}^*(\mathbf{u}; \phi)$ to the base $p_{\mathbf{u}}(\mathbf{u}; \psi)$ under the **reverse KL divergence**."

3-3-4. Alternative Divergences

not restricted to KL-divergence

- 1) f -divergence : use **density ratios** to compare models
- 2) Integral Probability Metrics (IPMs) : use **differences** for comparison

$$f\text{-divergence} \quad D_f [p_x^*(\mathbf{x}) \| p_x(\mathbf{x}; \boldsymbol{\theta})] = \mathbb{E}_{p_x(\mathbf{x}; \boldsymbol{\theta})} \left[f \left(\frac{p_x^*(\mathbf{x})}{p_x(\mathbf{x}; \boldsymbol{\theta})} \right) \right]$$

$$\text{IPM} \quad \delta_s [p_x^*(\mathbf{x}) \| p_x(\mathbf{x}; \boldsymbol{\theta})] = \mathbb{E}_{p_x^*(\mathbf{x})} [s(\mathbf{x})] - \mathbb{E}_{p_x(\mathbf{x}; \boldsymbol{\theta})} [s(\mathbf{x})]$$

4. Constructing Flows Part 1 : Finite Compositions

NF are **composable** : $T = T_K \circ \dots \circ T_1$

- idea : use simple transformation as building blocks

Evaluation

- forward evaluation : $\mathbf{z}_k = T_k(\mathbf{z}_{k-1})$ for $k = 1 : K$
- inverse evaluation : $\mathbf{z}_{k-1} = T_k^{-1}(\mathbf{z}_k)$ for $k = K : 1$

Jacobian-determinant (in log domain) : $\log |J_T(\mathbf{z})| = \log \left| \prod_{k=1}^K J_{T_k}(\mathbf{z}_{k-1}) \right| = \sum_{k=1}^K \log |J_{T_k}(\mathbf{z}_{k-1})|$

.

- increase depth \rightarrow only $O(K)$ growth in computational complexity

Implement T_k or T_k^{-1} using NN with param $\phi_k = f_{\phi_k}$

(ensure network is invertible & tractable Jacobian determinant)

Ensuring f_{ϕ_k} is invertible \neq Explicitly calculating its inverse

Tractable Jacobian determinant

- we can always find it but cost $O(D^3)$ (with D inputs & D outputs)
 \rightarrow intractable for large D
- should be at most $O(D)$
will describe NN design that allow Jacobian determinant to be computed in linear time

4-1. Autoregressive Flow

under certain conditions, we can transform **any distn** $p_x(\mathbf{x})$ into a **uniform distn** using maps with **triangular Jacobian**

Autoregressive Flow

$z'_i = \tau(z_i; h_i)$ where $h_i = c_i(\mathbf{z}_{<i})$.

- τ : transformer

- strictly monotonic function of z_i (thus invertible)
- parameterized by h_i
- specifies how the flow acts on z_i to output z'_i
- c_i : i -th conditioner
 - determines the parameters of the transformer
 - DOES NOT need to be a bijection

The above is **invertible** for any choice of τ and c_i , as long as "transformer is invertible"

$$z_i = \tau^{-1}(z'_i; \mathbf{h}_i) \quad \text{where} \quad \mathbf{h}_i = c_i(\mathbf{z}_{<i}).$$

- (forward : $z_1 \rightarrow z_K$) can be done in parallel
- (inverse : $z_k \rightarrow z_1$) all $\mathbf{z}_{<i}$ need to be computed before \mathbf{z}_i

Triangular Jacobian \rightarrow tractable ($\mathcal{O}(D)$ time)

$$J_{f_\phi}(\mathbf{z}) = \begin{bmatrix} \frac{\partial \tau}{\partial z_1}(z_1; \mathbf{h}_1) & & 0 \\ & \ddots & \\ \mathbf{L}(\mathbf{z}) & & \frac{\partial \tau}{\partial z_D}(z_D; \mathbf{h}_D) \end{bmatrix}.$$

$$\log|\det J_{f_\phi}(\mathbf{z})| = \log\left|\prod_{i=1}^D \frac{\partial \tau}{\partial z_i}(z_i; \mathbf{h}_i)\right| = \sum_{i=1}^D \log \left| \frac{\partial \tau}{\partial z_i}(z_i; \mathbf{h}_i) \right|$$

Autoregressive flows are **Universal approximators** , provided that the transformer & conditioner are flexible enough!

Alternative : conditioner c_i take in $\mathbf{z}'_{<i}$ instead of $\mathbf{z}_{<i}$

(mathematically equivalent)

4-1-1. Implementing the Transformer

What to choose as a transformer?

(a) Affine autoregressive flows

(= **location-scale** transformation)

$$\text{flow} : \tau(z_i; \mathbf{h}_i) = \alpha_i z_i + \beta_i \quad \text{where} \quad \mathbf{h}_i = \{\alpha_i, \beta_i\}.$$

$$\text{Jacobian} : \log|\det J_{f_\phi}(\mathbf{z})| = \sum_{i=1}^D \log|\alpha_i| = \sum_{i=1}^D \tilde{\alpha}_i.$$

Pros & Cons

- Pros) simplicity & analytical tractability
- Cons) expressivity is limited

(ex. let \mathbf{z} : Gaussian \rightarrow output is also GaussianThus need to stack multiple affine AF layers)

Widely used

- NICE, Real NVP, IAF, MAF, Glow

(b) Non-Affine neural transformers

conic combinations & compositions of monotonic functions are also monotonic

- Conic combination: $\tau(\mathbf{z}) = \sum_{k=1}^K w_k \tau_k(\mathbf{z})$, where $w_k > 0$ for all k .
- Composition: $\tau(\mathbf{z}) = \tau_K \circ \dots \circ \tau_1(\mathbf{z})$

Non-Affine neural transformers can be constructed...

"using **conic combination of monotonically increasing activation function** $\sigma(\cdot)$ "

$$\tau(\mathbf{z}_i; \mathbf{h}_i) = w_{i0} + \sum_{k=1}^K w_{ik} \sigma(\alpha_{ik} \mathbf{z}_i + \beta_{ik}) \quad \text{where} \quad \mathbf{h}_i = \{w_{i0}, \dots, w_{iK}, \alpha_{ik}, \beta_{ik}\}.$$

Jacobian determinant

- analytically obtainable (but more commonly computed by **back-propagation**)

Drawback

- (in general) **can not be inverted analytically**
(can be inverted only iteratively... e.g. **bijection search**)

Variants : **NAF, B-NAF, Flow++**

(c) Integration-based transformers

The integral of some positive function = monotonically increasing function

ex) $\tau(\mathbf{z}_i; \mathbf{h}_i) = \int_0^{\mathbf{z}_i} g(\mathbf{z}; \boldsymbol{\alpha}_i) d\mathbf{z} + \beta_i \quad \text{where} \quad \mathbf{h}_i = \{\boldsymbol{\alpha}_i, \beta_i\}.$

- $g(\cdot; \boldsymbol{\alpha}_i)$: any positive-valued NN
- derivative of transformer = $g(\mathbf{z}_i; \boldsymbol{\alpha}_i)$
- $g(\cdot; \boldsymbol{\alpha}_i)$ to be a positive polynomial of degree $2L$ (then integral will be $2L + 1$ in \mathbf{z}_i)
 \rightarrow **Sum-of-squares polynomial transformer**

$$\tau(\mathbf{z}_i; \mathbf{h}_i) = \int_0^{\mathbf{z}_i} \sum_{k=1}^K \left(\sum_{\ell=0}^L \alpha_{ik\ell} \mathbf{z}^\ell \right)^2 d\mathbf{z} + \beta_i.$$

Affine transformer : $L = 0$ of above

$$\int_0^{z_i} \sum_{k=1}^K (\alpha_{ik0} z^0)^2 dz + \beta_i = \left(\sum_{k=1}^K \alpha_{ik0}^2 \right) z \Big|_0^{z_i} + \beta_i = \alpha_i z_i + \beta_i.$$

(for large enough L) can approximate arbitrarily well any monotonic increasing function

(d) Neural spline flows

Non-affine transformers...don't have analytic inverse

But since all transformers are monotonic, can be inverted by **bijection search**

- $\tau(\mathbf{z}_{iA}; \mathbf{h}_i) < \mathbf{z}'_i < \tau(\mathbf{z}_{iB}; \mathbf{h}_i)$.
- keep halving, such that solution \mathbf{z}_i is always contained
- Bijection search computes \mathbf{z}_i with accuracy ϵ in $\mathcal{O}(\log \frac{1}{\epsilon})$
→ trade-off btw accuracy & computation

Overcome this trade-off with **monotonic spline**

(= piecewise function consisting of K segments , which are easy to invert)

- given $K + 1$ input locations $\mathbf{z}_{i0}, \dots, \mathbf{z}_{iK}$
- transformer $\tau(\mathbf{z}_i; \mathbf{h}_i)$: simple monotonic function in each interval $[\mathbf{z}_{i(k-1)}, \mathbf{z}_{ik}]$
- parameters \mathbf{h}_i :
 - input locations z_{i0}, \dots, z_{iK}
 - corresponding output locations z'_{i0}, \dots, z'_{iK}
 - derivatives (i.e. slopes) at $\mathbf{z}_{i0}, \dots, \mathbf{z}_{iK}$

Spline-based transformers

- distinguished by the type of spline they use
- fast to invert as to evaluate, while maintaining exact analytical invertibility
- Evaluating or inverting
 - (step 1) locate right segment ($\mathcal{O}(\log K)$ using binary search)
 - (step 2) evaluate / invert that segment (= analytically tractable)

4-1-2. Implementing the Conditioner

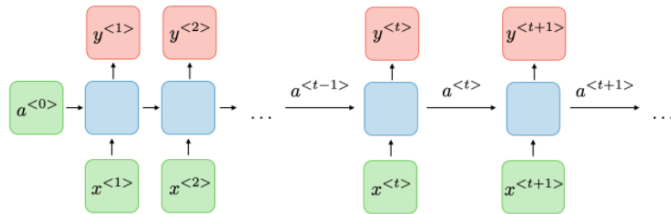
Conditioner : $c_i(\mathbf{z}_{<i})$

- can be any function (NN ok)
- each $c_i(\mathbf{z}_{<i})$ having its separate NN → scale poorly with dim D
(cost of storing, learning params of D independent networks...)
- should SHARE params! how?
 - 1) Recurrent Autoregressive Flows
 - 2) Masked Autoregressive Flows
 - 3) Coupling Layer

(a) Recurrent Autoregressive Flows

conditioner = RNN

$\mathbf{h}_i = c(\mathbf{s}_i)$ where $\mathbf{s}_1 = \text{initial state}$
 $\mathbf{s}_i = \text{RNN}(\mathbf{z}_{i-1}, \mathbf{s}_{i-1})$ for $i > 1$



share params across conditional distributions of **autoregressive models!**

Downside?

- parallel computation into **sequential** one
- involves $O(D)$ steps...slow for high-dim data

(b) Masked Autoregressive Flows

"Share params + avoid sequential computation of RNN by **MAKING** "

by Masking...

- single feedforward! (takes in \mathbf{z} , outputs the entire sequence $(\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_D)$)
- only requirement : **autoregressive structure**
(output \mathbf{h}_i can not depend on inputs $\mathbf{z}_{\geq i}$)
- remove connections from input z_i to outputs $(\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_i)$
(by multiplying each weight matrix elementwise, with binary matrix of same size)
- will have same architecture as original NN
- evaluate efficiently using GPU
- ex) **MADE** (Masked Autoencoder for Distribution Estimation), CNN, Self-attention...

MADE (Masked Autoencoder for Distribution Estimation)

- assign a "degree" between $1 \sim D$ for each input/hidden/output node
- mask-out the weights between subsequent layers

2 main advantages

- 1) Efficient to evaluate (one NN pass & can be computed parallel)
- 2) Universal approximators

Disadvantage

- Not as efficient to **invert** as to evaluate

(\because params \mathbf{h}_i are needed to obtain $\tau^{-1}(\mathbf{z}'_i; \mathbf{h}_i)$, thus can not be computed until (z_i, \dots, z_{i-1}) is obtained)

Pseudocode

Initialize z to an arbitrary value

for $i = 1 : D$

$(\mathbf{h}_1, \dots, \mathbf{h}_D) = c(\mathbf{z})$

$\mathbf{z}_i = \tau^{-1}(\mathbf{z}'_i; \mathbf{h}_i)$

- inverse : D times more expensive than evaluating forward transformation

Examples of masking to implementing **autoregressive flows**

- IAF, MAF, NAF, B-NAF, MintNet, MaCow

Examples of masking to implementing **non-flow based autoregressive models**

- MADE, PixelCNN, WaveNet

(c) Coupling Layers

Masked Autoregressive Flows : computational asymmetry

(either (1) sampling or (2) density evaluation will be D times slower IAF vs MAF)

→ for both to be fast, different conditioner is needed! ... **COUPLING LAYER**

Coupling layer

- choose index d (commonly, $D/2$)
- design the conditioner such that
 - params $(\mathbf{h}_1, \dots, \mathbf{h}_d)$ are constants
 - params $(\mathbf{h}_{d+1}, \dots, \mathbf{h}_D)$ are functions of $\mathbf{z}_{\leq d}$ (ex. NN)
- (coupling layer) splits into 2 parts
 - (fully autoregressive flow) splits into D parts
 - (intermediate) split into K parts

→ inverting the transformation will be $O(K)$ times more expensive than evaluating
- [forward]

$$\begin{aligned} \mathbf{z}'_{\leq d} &= \mathbf{z}_{\leq d} \\ (\mathbf{h}_{d+1}, \dots, \mathbf{h}_D) &= \text{NN}(\mathbf{z}_{\leq d}) \\ \mathbf{z}'_i &= \tau(\mathbf{z}_i; \mathbf{h}_i) \text{ for } i > d. \end{aligned}$$

- [inverse]

$$\begin{aligned} \mathbf{z}_{\leq d} &= \mathbf{z}'_{\leq d} \\ (\mathbf{h}_{d+1}, \dots, \mathbf{h}_D) &= \text{NN}(\mathbf{z}_{\leq d}) \\ \mathbf{z}_i &= \tau^{-1}(\mathbf{z}'_i; \mathbf{h}_i) \text{ for } i > d. \end{aligned}$$

- Jacobian :

$$J_{f_\phi} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{A} & \mathbf{D} \end{bmatrix}$$

- $\mathbf{I} : d \times d$ matrix
- $\mathbf{0} : d \times (D - d)$ matrix
- $\mathbf{A} : (D - d) \times d$ matrix
- $\mathbf{D} : (D - d) \times (D - d)$ matrix
- thus, jacobian : product of diagonals of $\mathbf{D} = \tau(\cdot; h_{d+1}), \dots, \tau(\cdot; h_D)$
- but due to efficiency...reduced expressive power
- single coupling layer \neq universal approximator
- need to compose multiple composing layers!
- (when composing, elements of \mathbf{z} need to be permuted! so that all dim have chance to be transformed)
- D coupling layers is an universal approximator
- (set index d of the i -th coupling layer = $i - 1$)
- Summary
 - coupling layers = most popular method for flow-based models
 - (\therefore allow both density evaluation & sampling to be done in single NN pass)
 - widely used in generative models of high-dim data
 - (NICE, Real NVP, Glow, WaveGlow, FloWaveNet, Flow++)

4-1-3. Relationship with Autoregressive Models

Alongside NF, **autoregressive models** are another popular model for high-dim distribution

$$p_{\mathbf{x}}(\mathbf{x}) = \prod_{i=1}^D p_{\mathbf{x}}(\mathbf{x}_i \mid \mathbf{x}_{<i}).$$

- model each conditional as $p_{\mathbf{x}}(\mathbf{x}_i \mid \mathbf{x}_{<i}) = p_{\mathbf{x}}(\mathbf{x}_i; \mathbf{h}_i)$, where $\mathbf{h}_i = c_i(\mathbf{x}_{<i})$
- ex) $p_{\mathbf{x}}(\mathbf{x}_i; \mathbf{h}_i)$: Gaussian, parameterized by its mean & var
- $c_i(\mathbf{x}_{<i})$: analogous to **conditioners** of autoregressive flow
- (= typically implemented with NN)
- can also be used for discrete/mixed data

Autoregressive models of continuous variables = **Autoregressive flows** with single layer

let $\tau(\mathbf{x}_i; \mathbf{h}_i)$ = cumulative distribution of $p_{\mathbf{x}}(\mathbf{x}_i; \mathbf{h}_i)$

$$\rightarrow \tau(\mathbf{x}_i; \mathbf{h}_i) = \int_{-\infty}^{\mathbf{x}_i} p_{\mathbf{x}}(\mathbf{x}'_i; \mathbf{h}_i) d\mathbf{x}'_i$$

and the vector $\mathbf{u} = (u_1, \dots, u_D)$

$$\rightarrow u_i = \tau(\mathbf{x}_i; \mathbf{h}_i) \quad \text{where} \quad \mathbf{h}_i = c_i(\mathbf{x}_{<i})$$

is always distributed uniformly in $(0, 1)^D$.

"same as autoregressive flow" with $\mathbf{z} = \mathbf{x}$ and $\mathbf{z}' = \mathbf{u}$.

Log probability :

$$\log p_{\mathbf{x}}(\mathbf{x}) = \log \prod_{i=1}^D \text{Uniform}(\tau(\mathbf{x}_i; \mathbf{h}_i); 0, 1) + \log \prod_{i=1}^D p_{\mathbf{x}}(\mathbf{x}_i; \mathbf{h}_i) = \sum_{i=1}^D \log p_{\mathbf{x}}(\mathbf{x}_i | \mathbf{x}_{<i}).$$

IAF : $\mathbf{z}_i = \tau^{-1}(\mathbf{u}_i; \mathbf{h}_i)$ where $\mathbf{h}_i = \mathbf{c}_i(\mathbf{x}_{<i})$

- corresponds exactly to sampling from autoregressive model (inverse transform sampling)

Transformer : not limited to being the inverse CDF

ex) $p_{\mathbf{x}}(\mathbf{x}_i; \mathbf{h}_i) = \mathcal{N}(\mathbf{x}_i; \mu_i, \sigma_i^2)$ where $\mathbf{h}_i = \{\mu_i, \sigma_i\}$

- reparam : $\mathbf{x}_i = \sigma_i \mathbf{u}_i + \mu_i$ where $\mathbf{u}_i \sim \mathcal{N}(0, 1)$
- thus, entire autoregressive model = affine autoregressive flow

Conclusion

- 1) AF = extending AM for continuous variables
- 2) Benefits of viewing AM as flows?
 - (a) provides a framework for their composition → enhance flexibility
 - (b) gives us freedom in specifying the base distn (base distn can be learned)
 - (c) compose AM with other types of flows (ex. non-autoregressive flows)

4-2. Linear Flows

Autoregressive flows : restrict \mathbf{z}'_i to depend only on inputs $\mathbf{z}_{\neq i}$

- dependent on the **order of input var**
- in the limit of capacity, no limit in flexibility, but in practice , it isn't
- So use **permutation!**

Permutation

- (1) easily invertible transformation
- (2) absolute Jacobian determinant = 1
- general idea of permutation = "Linear Flow"

$$\mathbf{z}' = \mathbf{W}\mathbf{z}$$

- where \mathbf{W} is a $D \times D$ invertible matrix that parameterizes the transformation
- Jacobian determinant = $\det \mathbf{W}$
- special case of linear flow, where \mathbf{W} is a **permutation matrix**
- Commonly... alternate "invertible linear transformation" & "autoregressive/coupling layers"

Parameterize and learn matrix \mathbf{W} ?

- problem 1) \mathbf{W} is not guaranteed to be invertible
- problem 2) solving $\mathbf{W}\mathbf{z} = \mathbf{z}'$ takes $O(D^3)$

problem 3) $\det \mathbf{W}$ takes $O(D^3)$ in general

- goal 1) \mathbf{W} that is invertible (& inversion cost $O(D^2)$)
- goal 2) computing its determinant costs $O(D)$

4-3. Residual Flows

invertible transformation of the form : $\mathbf{z}' = \mathbf{z} + g_\phi(\mathbf{z})$

- g_ϕ : NN that outputs D dim translation vector
- ϕ : NN parameter

2 general approaches

- 1) contractive maps
- 2) matrix determinant lemma

4-3-1. Contractive RF

RF is not always invertible, but invertible when g_ϕ **can be made contractive w.r.t some distance function**

- $\delta(F(\mathbf{z}_A), F(\mathbf{z}_B)) \leq L \delta(\mathbf{z}_A, \mathbf{z}_B)$.
- meaning) contractive map brings any 2 inputs close together, by at least a factor L

Banach fixed-point theorem

- contractive map has exactly one fixed point $\mathbf{z}_* = F(\mathbf{z}_*)$
- $\mathbf{z}_{k+1} = F(\mathbf{z}_k)$

Residual transformation : $\mathbf{z}' = f_\phi(\mathbf{z}) = \mathbf{z} + g_\phi(\mathbf{z})$

- $F(\hat{\mathbf{z}}) = \mathbf{z}' - g_\phi(\hat{\mathbf{z}})$
- If g_ϕ is contractive with Lipschitz constant L , then F is also \sim .
- by *Banach fixed-point theorem*,
 - there exists a unique \mathbf{z}_* such that $\mathbf{z}_* = \mathbf{z}' - g_\phi(\mathbf{z}_*)$
 - starting from arbitrary input \mathbf{z}_0 , iteratively apply F as follows:
$$\mathbf{z}_{k+1} = \mathbf{z}' - g_\phi(\mathbf{z}_k) \quad \text{for } k \geq 0$$
 - guarantees that the procedure above converges to $\mathbf{z}_* = f_\phi^{-1}(\mathbf{z}')$
 - rate of convergence : $\delta(\mathbf{z}_k, \mathbf{z}_*) \leq \frac{L^k}{1-L} \delta(\mathbf{z}_0, \mathbf{z}_1)$
smaller L , faster \mathbf{z}_k converges to \mathbf{z}_*
(trading off flexibility for efficiency)

Challenge of **contractive RF**

how to make contractive, without impinging upon its flexibility?

- composition of K Lipschitz-cont func F_1, \dots, F_K = Lipschitz cont with constant equal to $\prod_{k=1}^K L_k$
 - where L_k is the Lipschitz constant of F_k
 - make each layer Lipschitz continuous with $L_k \leq 1$ & at least one $L_k < 1$
 - ex) sigmoid, tanh, ReLU : Lipschitz continuous with a constant ≤ 1
- linear layer can be made **contractive**, by dividing them with operator norm
 - ex) Spectral normalization : use euclidean norm

Drawback of **contractive RF**

- no known general, efficient procedure for computing Jacobian costs $O(D^3)$
- nonetheless, able to obtain unbiased estimate of log absolute Jacobian
(by power series)

$$\log |\det J_{f_\phi}(\mathbf{z})| = \log \left| \det \left(\mathbf{I} + J_{g_\phi}(\mathbf{z}) \right) \right| = \sum_{k=1}^{\infty} \frac{(-1)^{k+1}}{k} \text{Tr} \{ J_{g_\phi}^k(\mathbf{z}) \}$$

- $J_{g_\phi}^k(\mathbf{z})$: k -th power of the Jacobian of g_ϕ evaluated at \mathbf{z}

(by Hutchinson trace estimator) efficiently estimate Trace

$$\text{Tr} \{ J_{g_\phi}^k(\mathbf{z}) \} \approx \mathbf{v}^\top J_{g_\phi}^k(\mathbf{z}) \mathbf{v}$$

- \mathbf{v} : D-dim random vector with zero mean & unit cov
- $\mathbf{v}^\top J_{g_\phi}^k(\mathbf{z})$ can be computed with k backprop

(by Russian-roulette estimator)

infinite sum can be estimated as finite sum of re-weighted terms

contractive RF (vs AF)

- dense Jacobian ... allows all input variables to affect all output variables
→ thus can be very flexible
- but exact density evaluation is computationally expensive &
sampling is done iteratively

4-3-2. RF based on the matrix determinant lemma

Matrix determinant lemma (= MD lemma)

$$\det(\mathbf{A} + \mathbf{V}\mathbf{W}^\top) = \det(\mathbf{I} + \mathbf{W}^\top \mathbf{A}^{-1} \mathbf{V}) \det \mathbf{A}$$

- left) $\mathcal{O}(D^3 + D^2 M)$
- right) $\mathcal{O}(M^3 + DM^2)$

(a) Planar Flow

g_ϕ is a one-layer NN with a single hidden unit:

$$\mathbf{z}' = \mathbf{z} + \mathbf{v}\sigma(\mathbf{w}^\top \mathbf{z} + b)$$

- $\mathbf{v} \in \mathbb{R}^D$, $\mathbf{w} \in \mathbb{R}^D$ and $b \in \mathbb{R}$, and σ is a differentiable activation function

Jacobian :

- $J_{f_\phi}(\mathbf{z}) = \mathbf{I} + \sigma'(\mathbf{w}^\top \mathbf{z} + b) \mathbf{v}\mathbf{w}^\top$.
- (using MD lemma) $\det J_{f_\phi}(\mathbf{z}) = 1 + \sigma'(\mathbf{w}^\top \mathbf{z} + b) \mathbf{w}^\top \mathbf{v}$

Conditions of Planar Flow to be invertible

- 1) σ' : positive & bounded above
- 2) $\mathbf{w}^\top \mathbf{v} > -\frac{1}{\sup_x \sigma'(x)}$
→ ensures $\det J_{f_\phi}(\mathbf{z})$ is always non-zero

(b) Sylvester Flow

extend Planar Flow to M hidden units

$$\mathbf{z}' = \mathbf{z} + \mathbf{V}\sigma(\mathbf{W}^\top \mathbf{z} + \mathbf{b}).$$

- $\mathbf{V} \in \mathbb{R}^{D \times M}$, $\mathbf{W} \in \mathbb{R}^{D \times M}$ and $\mathbf{b} \in \mathbb{R}^M$, and activation function σ to be elementwise

Jacobian

- $J_{f_\phi}(\mathbf{z}) = \mathbf{I} + \mathbf{V}\mathbf{S}(\mathbf{z})\mathbf{W}^\top$
- (by MD lemma) $\det J_{f_\phi}(\mathbf{z}) = \det(\mathbf{I} + \mathbf{S}(\mathbf{z})\mathbf{W}^\top \mathbf{V})$

QR decomposition

- $\mathbf{V} = \mathbf{Q}\mathbf{U}$ and $\mathbf{W} = \mathbf{Q}\mathbf{L}$
 - $\mathbf{Q} : D \times M$ matrix...whose columns are orthonormal ($\mathbf{Q}^T \mathbf{Q} = \mathbf{I}$)
 - $\mathbf{U} : M \times M$ upper triangle
 - $\mathbf{L} : M \times M$ lower triangle
- Thus, Jacobian determinant :

$$\det J_{f_\phi}(\mathbf{z}) = \det(\mathbf{I} + \mathbf{S}(\mathbf{z})\mathbf{L}^\top \mathbf{U}) = \prod_{i=1}^D (1 + S_{ii}(\mathbf{z})L_{ii}U_{ii})$$

(c) Radial flow

$$\mathbf{z}' = \mathbf{z} + \frac{\beta}{\alpha + r(\mathbf{z})}(\mathbf{z} - \mathbf{z}_0) \quad \text{where} \quad r(\mathbf{z}) = \|\mathbf{z} - \mathbf{z}_0\|$$

- $\alpha \in (0, +\infty)$, $\beta \in \mathbb{R}$ and $\mathbf{z}_0 \in \mathbb{R}^D$, and $\|\cdot\|$ is the norm
- (meaning) contraction/expansion radially with center \mathbf{z}_0

Jacobian :

- $J_{f_\phi}(\mathbf{z}) = \left(1 + \frac{\beta}{\alpha + r(\mathbf{z})}\right) \mathbf{I} - \frac{\beta}{r(\mathbf{z})(\alpha + r(\mathbf{z}))^2} (\mathbf{z} - \mathbf{z}_0) (\mathbf{z} - \mathbf{z}_0)^\top$.
- (using MD lemma) $\det J_{f_\phi}(\mathbf{z}) = \left(1 + \frac{\alpha\beta}{(\alpha + r(\mathbf{z}))^2}\right) \left(1 + \frac{\beta}{\alpha + r(\mathbf{z})}\right)^{D-1}$
be computed in $\mathcal{O}(D)$

The radial flow is not invertible for all β ...condition?

- " $\beta > -\alpha$ "ensures that $\det J_{f_\phi}(\mathbf{z})$ is always non-zero

Summary of planar, Sylvester, radial flows

- $\mathcal{O}(D)$ Jacobian determinant
- can be made invertible by suitably restricting their params
(but no analytical way to compute inverse)
- \therefore mostly used to **approximate posterior** for VAE (not in generative model)

4-4. Practical considerations when combining transformations

compose many flows as possible?

Glow

- 320 sub-transformations & 40 GPUs
- challenging computation....

2 techniques to solve?

- 1) stabilize the optimization
- 2) ease with computational demands

(a) Batch Normalization

stabilize & improve NN

composition of 2 affine transformation

- 1) scale & translation params (by batch statistics)
- 2) free params α (scale) & β (translation)

$$\text{BN}(\mathbf{z}) = \alpha \odot \frac{\mathbf{z} - \hat{\mu}}{\sqrt{\hat{\sigma}^2 + \epsilon}} + \beta, \quad \text{BN}^{-1}(\mathbf{z}') = \hat{\mu} + \frac{\mathbf{z}' - \beta}{\alpha} \odot \sqrt{\hat{\sigma}^2 + \epsilon}.$$

easy to compute Jacobian determinant (due to acting element-wise)

$$\det J_{\text{BN}}(\mathbf{z}) = \prod_{i=1}^D \frac{\alpha_i}{\sqrt{\hat{\sigma}_i^2 + \epsilon_i}}$$

Inserted between transformations! ($T_k \circ \text{BN} \circ T_{k-1}$)

ex) Glow

- employs a variant called **activation normalization**
- doesn't use batch statistics $\hat{\mu}$ and $\hat{\sigma}$.
(instead, before training, flow once & α and β are set so that batch has zero & unit var)
(= data-dependent initialization)
(α and β are optimized as model params)
- preferable when training with small-mini batches
(\therefore batch norm's statistics can be noisy)

(b) Multi-scale architecture

\mathbf{x} and \mathbf{u} should have same dim & T_k must preserve dimensionality

Multi-scale architecture :

- when going from \mathbf{x} to \mathbf{u} , some sub-dim of \mathbf{z}_k are clamped!
(& no additional transformation)
- just like **skip-connection**
- these help optimization!
- encode more global & semantic info in the dimensions

5. Constructing Flows part 2 : Continuous-Time Transformations

(until now) DISCRETE one-step transformation

(now) CONTINUOUS time, by parameterizing **flow's infinitesimal dynamics**

→ construct flow by defining **ODE** (Ordinary Differential Equation)
(describes flow's evolution over time...called "continuous time")

5-1. Definition

\mathbf{z}_t : flow's state at time t

- time t runs continuously from t_0 to t_1
- $\mathbf{z}_{t_0} = \mathbf{u}$ and $\mathbf{z}_{t_1} = \mathbf{x}$
- "parameterize the time derivative of \mathbf{z}_t " with NN g_ϕ

$$\frac{d\mathbf{z}_t}{dt} = g_\phi(t, \mathbf{z}_t).$$

$$\text{ODE} : \frac{d\mathbf{z}_t}{dt} = g_\phi(t, \mathbf{z}_t)$$

- requirement of g_ϕ : uniformly Lipschitz continuous in \mathbf{z}_t & continuous in t
(many NN meet those requirements)
- do not need invertibility & tractability of Jacobian determinant

to compute transformation ($\mathbf{x} = T(\mathbf{u})$)

- need to run dynamics forward in time
(= need to do integration)
- $\mathbf{x} = \mathbf{z}_{t_1} = \mathbf{u} + \int_{t=t_0}^{t_1} g_\phi(t, \mathbf{z}_t) dt.$

inverse transform T^{-1}

- $\mathbf{u} = \mathbf{z}_{t_0} = \mathbf{x} + \int_{t=t_1}^{t_0} g_\phi(t, \mathbf{z}_t) dt = \mathbf{x} - \int_{t=t_0}^{t_1} g_\phi(t, \mathbf{z}_t) dt$
- **Unlike discrete compositions (in before sections)**
continuous-time flows have the same computational complexity in each direction

change in log density

- $\frac{d \log p(\mathbf{z}_t)}{dt} = -\text{Tr} \left\{ J_{g_\phi(t, \cdot)}(\mathbf{z}_t) \right\}.$
 - $\text{Tr}\{\cdot\}$: trace operator
 - $J_{g_\phi(t, \cdot)}(\mathbf{z}_t)$: Jacobian of $g_\phi(t, \cdot)$ at \mathbf{z}_t
- but trace operator requires $O(D)$ back-prop!

use "Hutchinson's trace estimator"

$$\text{Tr} \left\{ J_{g_\phi(t, \cdot)}(\mathbf{z}_t) \right\} \approx \mathbf{v}^\top J_{g_\phi(t, \cdot)}(\mathbf{z}_t) \mathbf{v}$$

- \mathbf{v} : any D -dim random vector with zero mean & unit covariance
- $\mathbf{v}^\top J_{g_\phi(t, \cdot)}(\mathbf{z}_t)$: computed in a single back-prop pass

which D times more efficient!

Integrating the derivative of $\log p(\mathbf{z}_t)$

- $\log p_{\mathbf{x}}(\mathbf{x}) = \log p_{\mathbf{u}}(\mathbf{u}) - \int_{t=t_0}^{t_1} \text{Tr} \left\{ J_{g_\phi(t, \cdot)}(\mathbf{z}_t) \right\} dt$
- forward transform & log density can be done simultaneously

$$\begin{bmatrix} \mathbf{x} \\ \log p_{\mathbf{x}}(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} \mathbf{u} \\ \log p_{\mathbf{u}}(\mathbf{u}) \end{bmatrix} + \int_{t=t_0}^{t_1} \begin{bmatrix} g_\phi(t, \mathbf{z}_t) \\ -\text{Tr} \left\{ J_{g_\phi(t, \cdot)}(\mathbf{z}_t) \right\} \end{bmatrix} dt.$$

- computation is not analytically feasible for general $g_\phi \rightarrow$ use numerical integration

5-2. Solving & Optimizing continuous-time flow

(1) Euler's method

(2) Adjoint method

5-2-1. Euler's method and equivalence to RF

Simplest numerical technique

$$\mathbf{z}_{t+\epsilon} \approx f_\phi(\mathbf{z}_t) = \mathbf{z}_t + \epsilon g_\phi(t, \mathbf{z}_t)$$

- ϕ can be optimized with gradients computed via back-prop
- looks like $\mathbf{z}' = \mathbf{z} + g_\phi(\mathbf{z})$ (= contractive residual flow)

Taylor series expansion

log absolute Jacobian determinant of f_ϕ as follows:

$$\text{like } \log|\det J_{f_\phi}(\mathbf{z})| = \log|\det(\mathbf{I} + J_{g_\phi}(\mathbf{z}))| = \sum_{k=1}^{\infty} \frac{(-1)^{k+1}}{k} \text{Tr}\{J_{g_\phi}^k(\mathbf{z})\}$$

$$\log|\det J_{f_\phi}(\mathbf{z}_t)| = \sum_{k=1}^{\infty} \frac{(-1)^{k+1} \epsilon^k}{k} \text{Tr}\{J_{g_\phi(t, \cdot)}^k(\mathbf{z}_t)\} = \epsilon \text{Tr}\{J_{g_\phi(t, \cdot)}(\mathbf{z}_t)\} + \mathcal{O}(\epsilon^2).$$

- $\log p(\mathbf{z}_{t+\epsilon}) = \log p(\mathbf{z}_t) - \epsilon \text{Tr}\{J_{g_\phi(t, \cdot)}(\mathbf{z}_t)\} + \mathcal{O}(\epsilon^2)$
- $\frac{\log p(\mathbf{z}_{t+\epsilon}) - \log p(\mathbf{z}_t)}{\epsilon} = -\text{Tr}\{J_{g_\phi(t, \cdot)}(\mathbf{z}_t)\} + \mathcal{O}(\epsilon)$
- becomes $\frac{d \log p(\mathbf{z}_t)}{dt} = -\text{Tr}\{J_{g_\phi(t, \cdot)}(\mathbf{z}_t)\}$ as $\epsilon \rightarrow 0$

5-2-2. The Adjoint Method

target $\mathcal{L}(\mathbf{x}; \phi)$

gradient of $\partial \mathcal{L} / \partial \mathbf{z}_t$

$$\frac{d}{dt} \left(\frac{\partial \mathcal{L}}{\partial \mathbf{z}_t} \right) = - \left(\frac{\partial \mathcal{L}}{\partial \mathbf{z}_t} \right)^\top \frac{\partial g_\phi(t, \mathbf{z}_t)}{\partial \mathbf{z}_t}$$

- \mathbf{z}_t : flow's intermediate state

gradient w.r.t ϕ

$$\frac{\partial \mathcal{L}(\mathbf{x}; \phi)}{\partial \phi} = \int_{t=t_1}^{t_0} \frac{\partial \mathcal{L}}{\partial \mathbf{z}_t} \frac{\partial g_\phi(t, \mathbf{z}_t)}{\partial \phi} dt$$

6. Applications

2 primitive operations :

- (1) density calculation
- (2) sampling

application to

- 1) probabilistic
- 2) inference
- 3) supervised learning
- 4) reinforcement learning

6-1. Probabilistic Modeling

assume finite number of draws \mathbf{x} from unknown generative process : $p_{\mathbf{x}}^*(\mathbf{x})$

- goal) make a good approximation to $p_{\mathbf{x}}^*(\mathbf{x})$
- popular method : MLE (use forward KL div)

$$\begin{aligned} D_{\text{KL}} [p_{\mathbf{x}}^*(\mathbf{x}) \| p_{\mathbf{x}}(\mathbf{x}; \theta)] &= -\mathbb{E}_{p_{\mathbf{x}}^*(\mathbf{x})} [\log p_{\mathbf{x}}(\mathbf{x}; \theta)] + \text{const} \\ &\approx -\frac{1}{N} \sum_{n=1}^N \log p_{\mathbf{x}}(\mathbf{x}_n; \theta) + \text{const} \\ &= -\frac{1}{N} \sum_{n=1}^N \log p_{\mathbf{u}}(T^{-1}(\mathbf{x}_n; \phi); \psi) + \log |J_{T^{-1}}(\mathbf{x}_n; \phi)| + \text{const}. \end{aligned}$$

With the resulting model $p_{\mathbf{x}}(\mathbf{x}; \theta)$

- 1) density estimation
- 2) generation

(a) Density Estimation

use model to calculate densities

- 1) low-dim cases ...NF can represent skewed, multi-modal densities
- 2) use density function to perform one-class classification
- 3) detect rotations and corruptions of images
- 4) composition allows for better density estimation (MADE, Real NVP)

(b) Generation

sampling from model (use as it is sampled from $p_{\mathbf{x}}^*(\mathbf{x})$)

image, video, etc...

6-2. Inference

modeling params to infer **unknown quantities** within a model

most common setting : $\int \pi(\eta) d\eta$

summarize the use of flows for

- 1) sampling
- 2) variational inference
- 3) likelihood-free inference

6-2-1. Importance and Rejection sampling

Importance Sampling (IS)

- auxiliary distn $q(\eta)$

$$\int \pi(\eta) d\eta = \int q(\eta) \frac{\pi(\eta)}{q(\eta)} d\eta = \mathbb{E}_{q(\eta)} \left[\frac{\pi(\eta)}{q(\eta)} \right] \approx \frac{1}{S} \sum_{s=1}^S \frac{\pi(\hat{\eta}_s)}{q(\hat{\eta}_s)}.$$
 - $q(\eta)$ is a user-specified density function
 - $\hat{\eta}_s$ is a sample from $q(\eta)$
- choice of $q(\eta)$
 - if it doesn't contain $\pi(\eta)$'s, estimate becomes impractical
 - if it is too broad : inefficiency
- Importance of using expressive proposals $q(\eta)$!
 - **employ NFs**
- Requires both sampling & density evaluation to be tractable for many flows

2 alternatives to optimize flow's params

- 1) minimize KL div
- 2) minimize the variance of IS estimator

Rejection sampling (RS)

- aims to draw samples from $p(\eta) = \pi(\eta)/Z$
- use Real NVP to parameterize proposal distn
 - (since coupling layers allow for fast density evaluation & sampling)

6-2-2. Reparameterizing models for MCMC

NF can be integrated into MCMC

→ by using the flow to **reparam the target distn**

(can effectively smooth away pathologies, by allowing MCMC to be run on the simpler & better-behaved base density)

Metropolis-Hastings ratio to the reparameterized model

$$r(\hat{\mathbf{u}}_*; \hat{\mathbf{u}}_t) = \frac{p_u(\hat{\mathbf{u}}_*)}{p_u(\hat{\mathbf{u}}_t)} = \frac{\pi(T(\hat{\mathbf{u}}_*; \phi)) |\det J_T(\hat{\mathbf{u}}_*; \phi)|}{\pi(T(\hat{\mathbf{u}}_t; \phi)) |\det J_T(\hat{\mathbf{u}}_t; \phi)|}.$$

- $\hat{\mathbf{u}}_*$: proposed value
- $\hat{\mathbf{u}}_t$: current value

6-2-3. Variational Inference

use NF to fit **distn over latent variables**

"flows can usefully serve as posterior approximations"

Use a (trained) flow-based model

$$p(\eta \mid \mathbf{x}) \approx q(\eta; \phi) = q_{\mathbf{u}}(\mathbf{u}) |\det J_T(\mathbf{u}; \phi)|^{-1}$$

- $q_{\mathbf{u}}(\mathbf{u})$: base distn
- $T(\cdot; \phi)$: transformation

By maximizing ELBO

$$\begin{aligned} \log p(\mathbf{x}) &\geq \mathbb{E}_{q(\eta; \phi)} [\log p(\mathbf{x}, \eta)] - \mathbb{E}_{q(\eta; \phi)} [\log q(\eta; \phi)] \\ &= \mathbb{E}_{q_{\mathbf{u}}(\mathbf{u})} [\log p(\mathbf{x}, T(\mathbf{u}; \phi))] - \mathbb{E}_{q_{\mathbf{u}}(\mathbf{u})} [\log q_{\mathbf{u}}(\mathbf{u})] + \mathbb{E}_{q_{\mathbf{u}}(\mathbf{u})} [\log |\det J_T(\mathbf{u}; \phi)|] \\ &= \mathbb{E}_{q_{\mathbf{u}}(\mathbf{u})} [\log p(\mathbf{x}, T(\mathbf{u}; \phi))] + \mathbb{H}[q_{\mathbf{u}}(\mathbf{u})] + \mathbb{E}_{q_{\mathbf{u}}(\mathbf{u})} [\log |\det J_T(\mathbf{u}; \phi)|] \end{aligned}$$

- $\mathbb{H}[q_{\mathbf{u}}(\mathbf{u})]$: differential entropy (constant w.r.t ϕ)
- with MC estimation...

$$\mathbb{E}_{q_{\mathbf{u}}(\mathbf{u})} [\log p(\mathbf{x}, T(\mathbf{u}; \phi))] \approx \frac{1}{S} \sum_{s=1}^S \log p(\mathbf{x}, T(\hat{u}_s; \phi))$$

$$\mathbb{E}_{q_{\mathbf{u}}(\mathbf{u})} [\log |\det J_T(\mathbf{u}; \phi)|] \approx \frac{1}{S} \sum_{s=1}^S \log |\det J_T(\hat{u}_s; \phi)|$$