

Inverse Autoregressive Flow (IAF)

Yichen Ji

University of Chicago
jiyichen@uchicago.edu

October 26, 2023

Presentation Overview

① VI, NF & how they interplay

Motivation

$NF \cap VI$

② Inverse Autoregressive Flow

Formulation

Complexity in Sampling & Density Evaluation

Cheatcode: Implementation in TFP

Structure in Q

VAE consensus: use NNs for both encoder and decoder

Simple case: $q_{\phi}(z | x) = \mathcal{N}\left(z | \mu_{\phi}(x), \text{diag}\left(\sigma_{\phi}^2(x)\right)\right)$

whose distribution parameters are outputs of the encoding MLP

More flexible specification of approximate posterior distributions?

- Structured mean-field
- Mixture model
 - require evaluation of log-likelihood and its gradients for each mixture component per parameter update = expensive
- **Normalizing Flows**
 - integrates with VAE: Encoder outputs not just the parameters of the base distribution $q_0(z)$, but also the parameters of the NF (later)

Great NF review paper as reference: Papamakarios et al. (2019) & Kobyzev et al. (2020)

Core - Change of Variables

Idea: A mechanism to construct new families of distributions by choosing an initial density and then chaining together (in)finite number of parameterized, invertible, differentiable transformations:

$$\mathbf{z}_0 \sim q(\mathbf{z}_0 | \mathbf{x}), \quad \mathbf{z}_i = \mathbf{f}_i(\mathbf{z}_{i-1}, \mathbf{x}) \quad \forall i = 1 \dots D \quad (1)$$

Since invertible and differentiable transformations are composable:

$$\log q(\mathbf{z}_D | \mathbf{x}) = \log q(\mathbf{z}_0 | \mathbf{x}) - \sum_{i=1}^D \log \det \left| \frac{d\mathbf{z}_i}{d\mathbf{z}_{i-1}} \right| \quad (2)$$

Note that $\left| \frac{d\mathbf{z}_i}{d\mathbf{z}_{i-1}} \right|$ quantifies the relative change of volume (how much multiplication by the matrix expands or contracts space)

Nice conditions to have

in order to be practically useful

- transformation f must be invertible;
 - why crucial? Related to sampling & density estimation (come back later)
- both f and $g = f^{-1}$ must be differentiable;
 - diffeomorphism!
- easy to compute density
 - easy to compute the determinant of Jacobian
- easy to sample from
 - possible to have both?
 - typically only achieves one of them

Flow & Normalizing Flow

Where are the names coming from?

- Flow: Trajectory that a collection of samples from the base distribution $q_0(z)$ being gradually transformed by the sequence of transformations f_1, \dots, f_K
 - Sampling operation $x = z_K = f(z_0)$, $f = f_K \circ \dots \circ f_1$
- Normalizing Flow: The inverse flow through $f_K^{-1}, f_{K-1}^{-1}, \dots, f_1^{-1}$ takes a collection of samples from $q_K(z)$ and 'normalizes' them into a collection of samples from a prescribed base $q_0(z)$ e.g. Multivariate Normal
 - Density evaluation operation $p_x(x) = p_z(f^{-1}(x)) \left| \det J_{f^{-1}}(x) \right|$

(board)

Illustration

and what's not told in figure

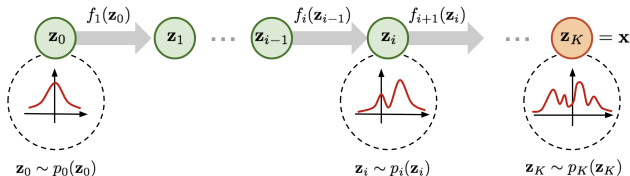


Figure: Generative direction - f pushes forward the base density $p_0(z_0)$ (noise)

The density of a sample can be evaluated by transforming it back to the base distribution and then computing the product of

- i) the density of the inverse-transformed sample under the base;
- ii) the associated change in volume induced by the sequence of inverse transformations.

Interplay

$NF \cap VI$

What VI looks like with NF? How NF interacts with VI for modeling and inference?

This point seems like a remaining question on which we haven't meditated enough.

I found Papamakarios et al. (2019) answered this question nicely, so I stole their ideas (and notations, sorry!).

Spoil alert: Reverse K-L divergence

NF \cap VI

with new notations...

Similarly to fitting any probabilistic model, fitting a flow-based model $p_x(x; \theta)$ to a target distribution $p_x^*(x)$ can be done by minimizing some divergence or discrepancy between them.

This minimization is performed with respect to the model's parameters $\theta = \{\phi, \psi\}$, where ϕ are the parameters of invertible transformation T and ψ are the parameters of base distribution $p_u(u)$.

think about VAE

In a standard Variational Autoencoder (VAE), the encoder neural network takes observed data x as input and outputs the parameters of a base distribution $p_u(u)$, commonly a Gaussian distribution. These parameters often include the mean and variance of the Gaussian.

After incorporating Normalizing Flows (NF) into a VAE, the encoder not only has to output the parameters of the base distribution $p_u(u)$, but also the parameters that define the transformations in the NF, i.e. $\theta = \{\phi, \psi\}$.

Forward K-L Divergence

Fancy name, but essentially 2 sides of K-L by asymmetry

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

Useful when: have samples from the target distribution (or the ability to generate them), but we cannot necessarily evaluate the target density $p_x^*(x)$

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

Reverse K-L

VI: reverse KL is all you need

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

Useful when: have the ability to evaluate the target density

$p_{\mathbf{x}}^*(x) = \frac{\tilde{p}_{\mathbf{x}}(x)}{C}$ up to a normalizing constant, but not necessarily sample from it.

$$\nabla_{\phi} \mathcal{L}(\theta) \approx -\frac{1}{N} \sum_{n=1}^N \nabla_{\phi} \log |\det J_T(\mathbf{u}_n; \phi)| + \nabla_{\phi} \log \tilde{p}_{\mathbf{x}}(T(\mathbf{u}_n; \phi))$$

Gradient estimate ∇_{ψ} by reparameterizing u as

$\mathbf{u} = T'(\mathbf{u}'; \psi)$ where $\mathbf{u}' \sim p_{\mathbf{u}'}(\mathbf{u}')$, like we did in VAE.

(Board)

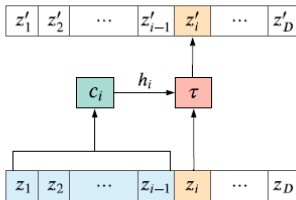
Inverse Autoregressive Flow

one example of flow construction

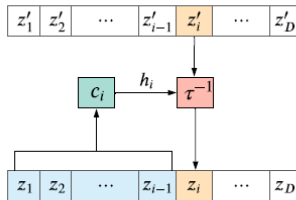
Autoregressive flow:

$$z'_i = \tau(z_i; \mathbf{h}_i) \quad \text{where} \quad \mathbf{h}_i = c_i(z_{<i})$$

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



(a) Forward



(b) Inverse

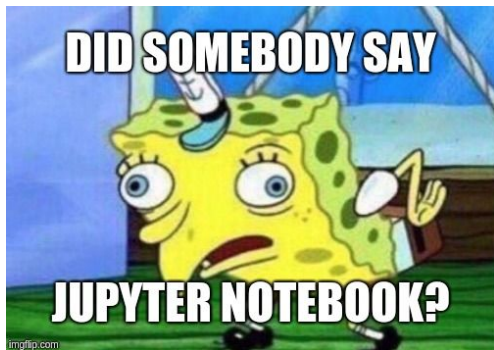
- Autoregressive Flow
 - Easy: Each z'_i can be computed in parallel (forward)
 - Hard: To compute z_i , all $z_{<i}$ need be have been computed (inverse)
- Inverse Autoregressive Flow
 - Easy: hard part of AF
 - Hard: easy part of AF

Puzzling inverse(to me, it's like inverse of inverse)!
(board to match notation in IAF paper)

Cheatcode

complexity in sampling and density evaluation

See Jupyter Notebook.



The End

Questions? Comments?