

Modern Approximate Inference: Variational Methods and Beyond



Diana Cai dcai@flatironinstitute.org



Yingzhen Li yingzhen.li@imperial.ac.uk

Why Probabilistic Inference?

Causal Inference and Causal Discovery



x: smokingm: tar in lungsy: lung cancer?

z: hidden confounder



How to estimate causal effects and/or discover causal relationships, when there exists hidden confounders?

Why Probabilistic Inference?

Uncertainty Quantification (Desiderata)

...

•••

...



Here's this patient's health record:

Could you suggest some next actions to improve her conditions?

Here's a list of potential treatment options based on the health record: [Treatment 1] with x% confidence (breakdown quantities)



Here's the conditions of this construction site:

Could you tell me what the potential safety issues are?

Here's a list of potential safety issues that need to be look after: [Issue 1] with x% confidence (breakdown quantities)





Why Probabilistic Inference?

Molecular Dynamics Simulations

- Energy of the molecular state: U(z)
- "Atoms never stop jiggling"
 - Sampling from the Boltzmann distribution (target):

$$\pi_T \propto \exp[-\frac{U(z)}{k_B T}]$$

• Computing "free energy"



Drug discovery

Material design



Energy (U)





The Central Computation Problem

- Inference: infer the unknowns
 - Unobserved/latent variables in the model
 - Quantities depending on the latent variables in the model



"statistics about the unknown"

• The central computational task for probabilistic inference:

$$\int F(z) \pi(z) dz$$

Causal Inference and Causal Discovery: Counterfactuals

"What is the counterfactual outcome y if I change the cause value to x' instead?

 $F(z) = p(y|x = x', z), \pi(z) = p(z | x, y),$ x, y = observed cause & outcome x' = alternative cause value



• The central computational task for probabilistic inference:

$$\int F(z) \pi(z) dz$$

Uncertainty Quantification

"What is the prediction distribution of the test output given a test input?"

 $F(z) = p(y|x, z), \pi(z) = p(z | D),$ D = observed datapoints



• The central computational task for probabilistic inference:

$$\int F(z) \pi(z) dz$$

Molecular Dynamics Simulations

"What is the free energy of a thermodynamic system at a particular end state?"

$$F(z) = 1, \pi(z) = \exp[-U(z)/k_BT],$$

$$\tilde{\pi} = \frac{1}{Z} \exp[-U(z)/k_BT]$$

$$-k_BT \log Z: \text{"free energy"}$$



(for biological conformational changes, ligand-macromolecule binding, chemical reaction mechanisms, etc.)

• The central computational task for probabilistic inference:

$$\int F(z) \pi(z) dz$$

Calculating Statistics

"What is the mean of this distribution?"

 $F(z) = z, \pi(z)$ can be complicated and high dimensional



• The central computational task for probabilistic inference:

$$\int F(z) \pi(z) dz$$

Evaluating the Likelihood

"What is the probability of generating this image?"

$$F(z) = N(net(z), \sigma^2 I), \pi(z) = N(0, I)$$



• The central computational task for probabilistic inference:

$$\int F(z) \pi(z) dz$$

Probabilistic Forecasting

"What is the weather forecast for tomorrow?"

Answering this in a Bayesian way: z: forecasting simulator settings D: historical weather record $F(z) = Simulator(z), \pi(z) = p(z | D)$



Nature laughs at the difficulties of integration.

-- Pierre-Simon Laplace

Gordon and Sorkin. The Armchair Science Reader. New York 1959



Integration in Bayesian Computation



Approximate Inference

• Central task: approximate $\pi(z)$

 $\overline{c} \qquad \overline{c}$ $q(z) \approx \pi(z)$



Approximate Inference

• Central task: approximate $\pi(z)$

Not simpler than computing $\int F(z) \pi(z) dz$ for a given integrand F(z)!

• Example: when $\pi(z) \propto \exp[-E(z)]$, estimating via self-normalized importance sampling

 $q(z) \approx \pi(z)$

$$\int F(z) \pi(z) dz = E_{\pi(z)}[F(z)] = E_{q(z)} \left[\frac{\pi(z)}{q(z)} F(z) \right] \approx \sum_{k=1}^{K} w_k F(z^k)$$
$$q(z) \gg \pi(z), z^k \sim q(z), \widetilde{w}_k \coloneqq \frac{\exp[-E(z^k)]}{q(z^k)}, w_k = \frac{\widetilde{w}_k}{\sum_{k'=1}^{K} \widetilde{w}_{k'}}$$

Solving $q(z) \approx \pi(z)$ can help computing $\int F(z) \pi(z) dz$ for any tractable integrand F(z):

$$\int F(z) \,\pi(z) dz = E_{\pi(z)}[F(z)] \approx E_{q(z)} \left[\frac{\pi(z)}{q(z)} F(z) \right] \approx \sum_{k=1}^{K} F(z^k), \, z^k \sim q(z)$$



Approximate Inference

• Central task: approximate $\pi(z)$



Approximate distribution design



Explicit distributions

SDE Trajectories Probability Flow

Distributions with samples but expensive/intractable density

The probability distribution
$$\pi(z)$$
 is intractable

$$\int_{a} \int_{a} \int_{a}$$

Algorithm for fitting q(z) to $\pi(z)$

min $Loss(q(z), \pi(z))$

Optimisation-based approaches

Sampling/Transport-based approaches

Previous Tutorial at NeurIPS 2020



Basics

Probabilistic modelling Approximate inference

Variational inference



Advances

Scalable variational inference

Monte Carlo techniques

Amortized inference

q distribution design

Optimization objective design



Applications

Bayesian neural networks Partially observed VAEs Future challenges



Prompt to Imagen 4

"A sci-fi image of two robots, one standing and doing calculations on whiteboard and another just sit in front of a computer, the computer screen shows codes, the whiteboard shows an integration equation $\int F(z) \pi(z) dz$ "



This Tutorial at UAI 2025 (5 Years Later)



Modern Basics

Variational inference Scaling & Monte Carlo Amortized inference Design principles I



New Advances

Score-based approximations Diffusion-based approximations Flow-based approximations Design principles II



Applications

Molecular Dynamics Simulation Simulation-based Inference Deep Generative Models Future challenges

Agenda for Today

- Basics
- (Coffee Break (30mins))
- Advances
- Small breaks (e.g., 5mins each) during sessions

Goal: Approximate an *intractable* target density $\pi(z)$ by a simpler family of densities \hat{Q}





Goal: Approximate an *intractable* target density $\pi(z)$ by a simpler family of densities \hat{Q}





Goal: Approximate an *intractable* target density $\pi(z)$ by a simpler family of densities \hat{Q}

Example: Bayesian inference

 $\pi(z) = p(z | \text{data}) = \frac{p(z) p(\text{data} | z)}{\int p(z) p(\text{data} | z) dz}$





Goal: Approximate an *intractable* target density $\pi(z)$ by a simpler family of densities \hat{Q}

Example: Bayesian inference

$$\pi(z) = p(z | data) = \frac{p(z) p(data | z)}{\int p(z) p(data | z)}$$

Often can't be computed in closed form







Goal: Approximate an *intractable* target density $\pi(z)$ by a simpler family of densities Q

Example: Bayesian inference

$$\pi(z) = p(z | \text{data}) = \frac{p(z) p(\text{data} | z)}{\int p(z) p(\text{data} | z)}$$

Often can't be computed in closed form

$$\pi(z) = \frac{\exp(-U(z)/(k_B T))}{\int \exp(-U(z)/(k_B T))dz}$$







Goal: Approximate an *intractable* target density $\pi(z)$ by a simpler family of densities Q

Example: Bayesian inference

$$\pi(z) = p(z | \text{data}) = \frac{p(z) p(\text{data} | z)}{\int p(z) p(\text{data} | z)}$$

Often can't be computed in closed form

$$\pi(z) = \frac{\exp(-U(z)/(k_B T))}{\int \exp(-U(z)/(k_B T))dz}$$









Goal: Approximate an *intractable* target density $\pi(z)$ by a simpler family of densities Q

Example: Bayesian inference

$$\pi(z) = p(z | \text{data}) = \frac{p(z) p(\text{data} | z)}{\int p(z) p(\text{data} | z)}$$

Often can't be computed in closed form

$$\pi(z) = \frac{\exp(-U(z)/(k_B T))}{\int \exp(-U(z)/(k_B T))dz}$$











Goal: Approximate an *intractable* target density $\pi(z)$ by a simpler family of densities Q

Example: Bayesian inference

$$\pi(z) = p(z | \text{data}) = \frac{p(z) p(\text{data} | z)}{\int p(z) p(\text{data} | z)}$$

Often can't be computed in closed form

$$\pi(z) = \frac{\exp(-U(z)/(k_B T))}{\int \exp(-U(z)/(k_B T))dz}$$













Goal: Approximate an *intractable* target density $\pi(z)$ by a simpler family of densities \hat{Q}





Goal: Approximate an *intractable* target density $\pi(z)$ by a simpler family of densities Q





Goal: Approximate an *intractable* target density $\pi(z)$ by a simpler family of densities Q

Example: Gaussian family with diagonal covariance $\mathcal{N}(\mu, \Sigma)$ (Often called "factorized" or "mean field" Gaussian)





Goal: Approximate an *intractable* target density $\pi(z)$ by a simpler family of densities Q

Example: Gaussian family with diagonal covariance $\mathcal{N}(\mu, \Sigma)$ (Often called "factorized" or "mean field" Gaussian)





[Zoltowski et al., 2021]





Goal: Approximate an *intractable* target density $\pi(z)$ by a simpler family of densities \hat{Q}





Goal: Approximate an *intractable* target density $\pi(z)$ by a simpler family of densities Q





Goal: Approximate an *intractable* target density $\pi(z)$ by a simpler family of densities \hat{Q}

Step 1: Specify a *divergence* $\mathscr{D}(q; \pi)$ between variational density q(z) and target $\pi(z)$







- **Step 1:** Specify a *divergence* $\mathscr{D}(q; \pi)$ between variational density q(z) and target $\pi(z)$
- **Step 2:** Find the member of the family \hat{Q} that minimizes the divergence, i.e.,

Goal: Approximate an *intractable* target density $\pi(z)$ by a simpler family of densities \hat{Q}






- **Step 1:** Specify a *divergence* $\mathscr{D}(q; \pi)$ between variational density q(z) and target $\pi(z)$
- **Step 2:** Find the member of the family \hat{Q} that minimizes the divergence, i.e.,

$$q^* = \arg\min_{q \in \mathcal{Q}} \mathcal{D}(q; \pi)$$

Goal: Approximate an *intractable* target density $\pi(z)$ by a simpler family of densities \hat{Q}







Goal: Approximate an *intractable* target density $\pi(z)$ by a simpler family of densities Q

Step 2: Find the member of the family \hat{Q} that minimizes the divergence, i.e.,

 $q^* = \arg\min_{q \in \mathcal{Q}} \mathcal{D}(q; \pi)$

Example: find the mean and (diagonal) covariance such that it is closest to π

Step 1: Specify a *divergence* $\mathscr{D}(q; \pi)$ between variational density q(z) and target $\pi(z)$

 $\mu^*, \Sigma^* = \arg\min \mathcal{D}(N(\mu, \Sigma); \pi)$ μ,Σ







Goal: Approximate an *intractable* target density $\pi(z)$ by a simpler family of densities Q

Step 2: Find the member of the family \hat{Q} that minimizes the divergence, i.e.,

 $q^* = \arg\min_{q \in \mathcal{Q}} \mathcal{D}(q; \pi)$

Example: find the mean and (diagonal) covariance such that it is closest to π

Step 1: Specify a *divergence* $\mathscr{D}(q;\pi)$ between variational density q(z) and target $\pi(z)$

 $\mu^*, \Sigma^* = \arg\min \mathcal{D}(N(\mu, \Sigma); \pi)$ μ,Σ







Goal: Approximate an *intractable* target density $\pi(z)$ by a simpler family of densities Q

Step 2: Find the member of the family \hat{Q} that minimizes the divergence, i.e.,

 $q^* = \arg\min_{q \in \mathcal{Q}} \mathcal{D}(q; \pi)$

Step 3: Use variational density q^* instead of target $\pi(z)$ in downstream tasks: e.g., $\mathbb{E}_{\pi(z)}[f(z)] \approx \mathbb{E}_{q^*(z)}[f(z)]$

Step 1: Specify a *divergence* $\mathscr{D}(q; \pi)$ between variational density q(z) and target $\pi(z)$

Example: find the mean and (diagonal) covariance such that it is closest to π

 $\mu^*, \Sigma^* = \arg\min \mathcal{D}(N(\mu, \Sigma); \pi)$ μ, Σ







Kullback-Leibler (KL) divergence ("forward KL")





Kullback-Leibler (KL) divergence ("forward KL")

 $\mathcal{D}_{KL}(\pi;q) \neq \mathcal{D}_{KL}(q;\pi)$







Kullback-Leibler (KL) divergence ("forward KL")



 $\mathscr{D}_{KL}(\pi;q) \neq \mathscr{D}_{KL}(q;\pi)$





Kullback-Leibler (KL) divergence ("forward KL")



$\mathscr{D}_{KL}(\pi;q) \neq \mathscr{D}_{KL}(q;\pi)$





Kullback-Leibler (KL) divergence ("forward KL")



 $\mathscr{D}_{KL}(\pi;q) \neq \mathscr{D}_{KL}(q;\pi)$ $\mathscr{D}_{KL}(\pi;q) = 0 \text{ iff } \pi = q$ **Problem:** we often don't know $\pi(z)$ - can't sample from it or evaluate it pointwise





Kullback-Leibler (KL) divergence ("forward KL")

$$\mathscr{D}_{KL}(\pi;q) = \int \log q$$

Most commonly used in VI: "reverse" KL divergence



 $\mathscr{D}_{KL}(\pi;q) \neq \mathscr{D}_{KL}(q;\pi)$ $g\left(\begin{array}{c}\pi(z)\\a(z)\end{array}\right)\pi(z)dz$ $\mathscr{D}_{KL}(\pi;q) = 0 \text{ iff } \pi = q$ **Problem:** we often don't Problem: we often don't know $\pi(z)$ - can't sample from it or evaluate it pointwise





Kullback-Leibler (KL) divergence ("forward KL")

$$\mathscr{D}_{KL}(\pi;q) = \int \log q$$

Most commonly used in VI: "reverse" KL divergence

$$\mathcal{D}_{KL}(q;\pi) = \int \log q$$



 $\mathscr{D}_{KL}(\pi;q) \neq \mathscr{D}_{KL}(q;\pi)$ $g\left(\frac{\pi(z)}{q(z)}\right)\pi(z)dz \qquad \qquad \mathscr{D}_{KL}(\pi;q) = 0 \text{ iff } \pi = q$ **Problem:** we often don't know $\pi(z)$ - can't sample from it or evaluate it pointwise

 $\log\left(\frac{q(z)}{\pi(z)}\right) q(z) dz$





Kullback-Leibler (KL) divergence ("forward KL")

$$\mathscr{D}_{KL}(\pi;q) = \int \log q$$

Most commonly used in VI: "reverse" KL divergence

$$\mathcal{D}_{KL}(q;\pi) = \int \log q$$



 $\mathscr{D}_{KL}(\pi;q) \neq \mathscr{D}_{KL}(q;\pi)$ $g\left(\begin{array}{c}\pi(z)\\q(z)\end{array}\right)\pi(z)dz\qquad \qquad \mathscr{D}_{KL}(\pi;q)=0 \text{ iff }\pi=q$ **Problem:** we often don't know $\pi(z)$ - can't sample from it or evaluate it pointwise



 $\int \log \left(\frac{q(z)}{\pi(z)} \right) q(z) dz$ But I still can't evaluate $\pi(z)$.





Kullback-Leibler (KL) divergence ("forward KL")

$$\mathscr{D}_{KL}(\pi;q) = \int \log q$$

Most commonly used in VI: "reverse" KL divergence

$$\mathcal{D}_{KL}(q;\pi) = \int \log q$$
$$= \int \log q$$



 $\mathscr{D}_{KI}(\pi;q) \neq \mathscr{D}_{KI}(q;\pi)$ $g\left(\frac{\pi(z)}{q(z)}\right)\pi(z)dz \qquad \qquad \mathscr{D}_{KL}(\pi;q) = 0 \text{ iff } \pi = q$ **Problem:** we often don't know $\pi(z)$ - can't sample from it or evaluate it pointwise



 $g\left(\frac{q(z)}{\pi(z)}\right) q(z) dz$ But I still can't evaluate $\pi(z)$.

 $g\left(\frac{q(z)}{\rho(z)}\right) q(z) dz + \text{constant}$





Kullback-Leibler (KL) divergence ("forward KL")

$$\mathscr{D}_{KL}(\pi;q) = \int \log q$$

Most commonly used in VI: "reverse" KL divergence

$$\mathcal{D}_{KL}(q;\pi) = \int \log q q dq$$
$$= \int \log q dq$$



 $\mathscr{D}_{KL}(\pi;q) \neq \mathscr{D}_{KL}(q;\pi)$ $g\left(\frac{\pi(z)}{q(z)}\right)\pi(z)dz \qquad \mathscr{D}_{KL}(\pi;q) = 0 \text{ iff } \pi = q$ **Problem: we often don't** know $\pi(z)$ - can't sample from it or evaluate it pointwise







Kullback-Leibler (KL) divergence ("forward KL")

$$\mathscr{D}_{KL}(\pi;q) = \int \log q$$

Most commonly used in VI: "reverse" KL divergence

$$\mathcal{D}_{KL}(q;\pi) = \int \log q$$



 $\mathscr{D}_{KL}(\pi;q) \neq \mathscr{D}_{KL}(q;\pi)$ $g\left(\begin{array}{c}\pi(z)\\a(z)\end{array}\right)\pi(z)dz$ $\mathscr{D}_{KL}(\pi;q) = 0 \text{ iff } \pi = q$ **Problem:** we often don't know $\pi(z)$ - can't sample from it or evaluate it pointwise



Quantities inside the second integral are **tractable**: I can evaluate $\rho(z)$, e.g., in Bayesian inference it's the joint distribution $\rho(z) = p(z) p(\text{data} | z)$





Forward KL $\mathscr{D}_{KL}(\pi; q)$



Murphy. Probabilistic Machine Learning. Book 1.

Reverse KL



























Here q "covers" π









Here q "covers" π

Here *q* latches onto one of the modes of π







The reverse KL divergence is intractable. Recall that we can write it in this tractable form:

 $\mathscr{D}_{KL}(q;p) = \left[\log\left(\frac{q(z)}{\rho(z)}\right) q(z) dz + \text{constant}\right]$





The reverse KL divergence is intractable. Recall that we can write it in this tractable form:

 $\mathcal{D}_{KL}(q;p) = \int \log \begin{pmatrix} q(z) \\ \rho(z) \end{pmatrix} q(z) dz + \text{constant}$ Bayes: $\rho(z) = p(z) p(\text{data} | z)$





The reverse KL divergence is intractable. Recall that we can write it in this tractable form:

$$\mathcal{D}_{KL}(q;p) = \int \log$$

 $g\begin{pmatrix} q(z) \\ \rho(z) \end{pmatrix} q(z) dz + \text{constant}$ Bayes: $\rho(z) = p(z) p(\text{data} | z)$ The *evidence lower bound* (ELBO) is the negative KL (up to an additive constant):





The reverse KL divergence is intractable. Recall that we can write it in this tractable form:

$$\mathscr{D}_{KL}(q;p) = \int \log \begin{pmatrix} q(z) \\ \rho(z) \end{pmatrix} q(z) dz + \text{constant}$$

Bayes: $\rho(z) = p(z) p(\text{data} | z)$

The evidence lower bound (ELBO) is the negative KL (up to an additive constant): $\mathsf{ELBO}(q) = \log\left(\frac{\rho(z)}{q(z)}\right) q(z) dz$





The reverse KL divergence is intractable. Recall that we can write it in this tractable form:

$$\mathscr{D}_{KL}(q;p) = \int \log \begin{pmatrix} q(z) \\ \rho(z) \end{pmatrix} q(z) dz + \text{constant}$$

Bayes: $\rho(z) = p(z) p(\text{data} | z)$

The evidence lower bound (ELBO) is the negative KL (up to an additive constant): $\mathsf{ELBO}(q) = \left[\log\left(\frac{\rho(z)}{q(z)}\right) q(z) dz\right]$

In practice, we maximize the ELBO:





The reverse KL divergence is intractable. Recall that we can write it in this tractable form:

$$\mathscr{D}_{KL}(q;p) = \int \log \begin{pmatrix} q(z) \\ \rho(z) \end{pmatrix} q(z) dz + \text{constant}$$

Bayes: $\rho(z) = p(z) p(\text{data} | z)$

The evidence lower bound (ELBO) is the negative KL (up to an additive constant): $\mathsf{ELBO}(q) = \left[\log\left(\frac{\rho(z)}{q(z)}\right) q(z) dz\right]$

In practice, we maximize the ELBO:

• max ELBO(q) = min $\mathscr{D}_{KL}(q;p)$ w.r.t the parameters of q





The reverse KL divergence is intractable. Recall that we can write it in this tractable form:

$$\mathscr{D}_{KL}(q;p) = \int \log \begin{pmatrix} q(z) \\ \rho(z) \end{pmatrix} q(z) dz + \text{constant}$$

Bayes: $\rho(z) = p(z) p(\text{data} | z)$

The evidence lower bound (ELBO) is the negative KL (up to an additive constant): $\mathsf{ELBO}(q) = \left[\log\left(\frac{\rho(z)}{q(z)}\right) q(z) dz\right]$

In practice, we maximize the ELBO:

- max ELBO(q) = min $\mathscr{D}_{KL}(q;p)$ w.r.t the parameters of q
- Not a convex objective





The reverse KL divergence is intractable. Recall that we can write it in this tractable form:

$$\mathcal{D}_{KL}(q;p) = \int \log \begin{pmatrix} q(z) \\ \rho(z) \end{pmatrix} q(z) dz + \text{constant}$$

Bayes: $\rho(z) = p(z) p(\text{data} | z)$

The evidence lower bound (ELBO) is the negative KL (up to an additive constant): $\mathsf{ELBO}(q) = \left[\log\left(\frac{\rho(z)}{q(z)}\right) q(z) dz\right]$

In practice, we maximize the ELBO:

- max ELBO(q) = min $\mathscr{D}_{KL}(q;p)$ w.r.t the parameters of q
- Not a convex objective

• Bayes: it is a *lower bound* on the marginal likelihood, i.e., $ELBO(q) \le \log p(data)$





The reverse KL divergence is intractable. Recall that we can write it in this tractable form:

$$\mathcal{D}_{KL}(q;p) = \int \log \left(\frac{q(z)}{\rho(z)} \right) q(z) dz + \text{constant}$$

Bayes: $\rho(z) = p(z) p(\text{data} | z)$

The evidence lower bound (ELBO) is the negative KL (up to an additive constant): $\mathsf{ELBO}(q) = \left[\log\left(\frac{\rho(z)}{q(z)}\right) q(z) dz = \log p(\mathsf{data}) - \mathscr{D}_{KL}(q(z); p(z | \mathsf{data}))\right]$

In practice, we maximize the ELBO:

- max ELBO(q) = min $\mathscr{D}_{KI}(q;p)$ w.r.t the parameters of q
- Not a convex objective

• Bayes: it is a *lower bound* on the marginal likelihood, i.e., $ELBO(q) \le \log p(data)$



26

The reverse KL divergence is intractable. Recall that we can write it in this tractable form:

$$\mathscr{D}_{KL}(q;p) = \int \log \begin{pmatrix} q(z) \\ \rho(z) \end{pmatrix} q(z) dz + \text{constant}$$

Bayes: $\rho(z) = p(z) p(\text{data} | z)$

The evidence lower bound (ELBO) is the negative KL (up to an additive constant):

$$\mathsf{ELBO}(q) = \int \log\left(\frac{\rho(z)}{q(z)}\right) q(z)$$

In practice, we maximize the ELBO:

- max ELBO(q) = min $\mathscr{D}_{KI}(q;p)$ w.r.t the parameters of q
- Not a convex objective

• Bayes: it is a *lower bound* on the marginal likelihood, i.e., ELBO(q) $\leq \log p(\text{data})$ Classical VI uses coordinate-ascent to maximize the ELBO

Blei et al. Variational Inference: A Review for Statisticians. Journal of the American Statistician, 2017.

 $dz = \log p(\text{data}) - \mathcal{D}_{KL}(q(z); p(z | \text{data}))$



26

Variational approximations with Gaussians





Many other divergences have been considered in VI

Some other divergences / distances considered in VI:

- α -divergence
- χ^2 -divergence
- Hellinger distance
- Kernelized Stein discrepancy
- Fisher divergence



Long history behind variational inference

Early work on approximate inference, e.g., Metropolis-Hastings, importance sampling



pre-90s





Jordan, Ghahramani, Jaakkola, Saul. Introduction to variational methods for graphical models. *Machine Learning*, 1999.





Long history behind variational inference



Jordan, Ghahramani, Jaakkola, Saul. Introduction to variational methods for graphical models. *Machine Learning*, 1999.

Blei, Ng, Jordan. Latent Dirichlet allocation. Journal of Machine Learning Research, 2003.





Long history behind variational inference



Jordan, Ghahramani, Jaakkola, Saul. Introduction to variational methods for graphical models. *Machine Learning*, 1999.

Blei, Ng, Jordan. Latent Dirichlet allocation. Journal of Machine Learning Research, 2003.

Hoffman, Blei, Wang, Paisley. Stochastic variational inference. Journal of Machine Learning Research, 2013. Ranganath, Gerrish, Blei. Black box variational inference. AISTATS, 2014. Rezende and Mohamed. Variational inference with normalizing flows. ICML, 2016. Kucukelbir, Tran, Ranganath, Gelman, Blei. Automatic differentiation variational inference. JMLR, 2018.




Long history behind variational inference



Jordan, Ghahramani, Jaakkola, Saul. Introduction to variational methods for graphical models. Machine Learning, 1999.

Blei, Ng, Jordan. Latent Dirichlet allocation. Journal of Machine Learning Research, 2003.

Hoffman, Blei, Wang, Paisley. Stochastic variational inference. Journal of Machine Learning Research, 2013. Ranganath, Gerrish, Blei. Black box variational inference. AISTATS, 2014. Rezende and Mohamed. Variational inference with normalizing flows. ICML, 2016. Kucukelbir, Tran, Ranganath, Gelman, Blei. Automatic differentiation variational inference. JMLR, 2018.





Long history behind variational inference

- Factorized (or "mean field") assumptions



Jordan, Ghahramani, Jaakkola, Saul. Introduction to variational methods for graphical models. *Machine Learning*, 1999.

Blei, Ng, Jordan. Latent Dirichlet allocation. Journal of Machine Learning Research, 2003.

Hoffman, Blei, Wang, Paisley. Stochastic variational inference. Journal of Machine Learning Research, 2013. Ranganath, Gerrish, Blei. Black box variational inference. AISTATS, 2014. Rezende and Mohamed. Variational inference with normalizing flows. ICML, 2016. Kucukelbir, Tran, Ranganath, Gelman, Blei. Automatic differentiation variational inference. JMLR, 2018.





Long history behind variational inference

Elements of "classical" VI:

- Factorized (or "mean field") assumptions
- Coordinate-ascent and ELBO maximization
- Restrictive assumptions on target (long derivations)



Rezende and Mohamed. Variational inference with normalizing flows. ICML, 2016. Kucukelbir, Tran, Ranganath, Gelman, Blei. Automatic differentiation variational inference. JMLR, 2018.



s 2010s

2020s

This tutorial: "modern" variational inference

"Black box" inference: few assumptions on target p(z)
Automatic differentiation; supported in modern software
Towards "expressive" variational families
Alternative divergences





Part II: Key components of modern methods



Towards black-box approaches for VI

Classical approaches to VI used coordinate-ascent VI (CAVI) updates.

- The variational family is factorized (a "mean field" family)
- Structural assumptions about the target, e.g., "conditionally conjugate" Challenging for practitioners to apply due to need for specialized derivations.

- In order to make the problem tractable, typically additional assumptions imposed:



Towards black-box approaches for VI

Classical approaches to VI used coordinate-ascent VI (CAVI) updates.

- The variational family is factorized (a "mean field" family)
- Structural assumptions about the target, e.g., "conditionally conjugate" Challenging for practitioners to apply due to need for specialized derivations.

- "Black-box" VI approaches: very few assumptions needed about the target • Bayesian modeling: usually just want the joint p(z, x) to be differentiable
 - We will see algorithms that use $\nabla_z \log p(z | x)$ ("the score")
- Autodiff tools: led to more user-friendly strategies for applying VI
- Implemented in many software packages

- In order to make the problem tractable, typically additional assumptions imposed:



Variational Bayes: learn an approximation for $\pi(z) = p(z | x)$ for data x





Variational Bayes: learn an approximation for $\pi(z) = p(z | x)$ for data x

Gradient-based variational inference:





Variational Bayes: learn an approximation for $\pi(z) = p(z | x)$ for data x

Gradient-based variational inference:

Variational parameters ϕ





Variational Bayes: learn an approximation for $\pi(z) = p(z | x)$ for data x

Gradient-based variational inference:

Variational parameters ϕ ELBC

$$D(\phi) = \int \log\left(\frac{p_{\theta}(z, x)}{q_{\phi}(z)}\right) q_{\phi}(z) dz$$





Variational Bayes: learn an approximation for $\pi(z) = p(z | x)$ for data x

Gradient-based variational inference:

ELBC Variational parameters ϕ

latent variables z.

$$D(\phi) = \int \log \left(\frac{p_{\theta}(z, x)}{q_{\phi}(z)} \right) q_{\phi}(z) dz$$

Joint distribution between *z* and *x*





Variational Bayes: learn an approximation for $\pi(z) = p(z | x)$ for data x

Gradient-based variational inference:

Variational parameters ϕ ELBO $(\phi) = \int \log \left(\frac{p_{\theta}(z, x)}{q_{\phi}(z)} \right) q_{\phi}(z) dz$ Will need gradients of the ELBO: ∇_{ϕ} ELBO(ϕ) Joint distribution between *z* and *x*





Variational Bayes: learn an approximation for $\pi(z) = p(z | x)$ for data x

Gradient-based variational inference:

Variational parameters ϕ ELBO $(\phi) = \int \log \left(\frac{p_{\theta}(z, x)}{q_{\phi}(z)} \right) q_{\phi}(z) dz$ Will need gradients of the ELBO: ∇_{ϕ} ELBO(ϕ) Joint distribution between *z* and *x*

What we will cover:





Variational Bayes: learn an approximation for $\pi(z) = p(z | x)$ for data x

Gradient-based variational inference:

ELBC Variational parameters ϕ

Will need gradients of the ELBO

What we will cover:

Monte Carlo approximations of the gradient

latent variables z.

$$D(\phi) = \int \log \left(\frac{p_{\theta}(z, x)}{q_{\phi}(z)} \right) q_{\phi}(z) dz$$

: $\nabla_{\phi} \text{ELBO}(\phi)$
Joint distribution between *z* and



X



Gradient-based variational inference:

ELBC Variational parameters ϕ

Will need gradients of the ELBO:

What we will cover:

Monte Carlo approximations of the gradient

Variational Bayes: learn an approximation for $\pi(z) = p(z | x)$ for data x

latent variables *z*.

$$D(\phi) = \int \log \left(\frac{p_{\theta}(z, x)}{q_{\phi}(z)} \right) q_{\phi}(z) dz$$

: $\nabla_{\phi} \text{ELBO}(\phi)$
Joint distribution between *z* and *x*

1) Monte Carlo approximation of ∇_{ϕ} ELBO(ϕ) via samples $z^1, \ldots, z^B \sim q_{\phi}$





Variational Bayes: learn an approximation for $\pi(z) = p(z | x)$ for data x

Gradient-based variational inference:

Variational parameters ϕ ELBC

Will need gradients of the ELBO:

What we will cover:

Monte Carlo approximations of the gradient

latent variables *z*.

$$D(\phi) = \int \log \left(\frac{p_{\theta}(z, x)}{q_{\phi}(z)} \right) q_{\phi}(z) dz$$

: $\nabla_{\phi} \text{ELBO}(\phi)$
Joint distribution between *z* and *x*

1) Monte Carlo approximation of ∇_{ϕ} ELBO(ϕ) via samples $z^1, \dots, z^B \sim q_{\phi}$





Variational Bayes: learn an approximation for $\pi(z) = p(z | x)$ for data x

Gradient-based variational inference:

ELBC Variational parameters ϕ

Will need gradients of the ELBO

What we will cover:

Monte Carlo approximations of the gradient

2) Score function gradient estimator

latent variables *z*.

$$D(\phi) = \int \log \left(\frac{p_{\theta}(z, x)}{q_{\phi}(z)} \right) q_{\phi}(z) dz$$

: $\nabla_{\phi} \text{ELBO}(\phi)$
Joint distribution between *z* and *x*

1) Monte Carlo approximation of ∇_{ϕ} ELBO(ϕ) via samples $z^1, \dots, z^B \sim q_{\phi}$





Variational Bayes: learn an approximation for $\pi(z) = p(z | x)$ for data x

Gradient-based variational inference:

ELBC Variational parameters ϕ

Will need gradients of the ELBO:

What we will cover:

Monte Carlo approximations of the gradient

- 2) Score function gradient estimator
- 3) Reparameterization gradient estimator

latent variables *z*.

$$D(\phi) = \int \log \left(\frac{p_{\theta}(z, x)}{q_{\phi}(z)} \right) q_{\phi}(z) dz$$

: $\nabla_{\phi} \text{ELBO}(\phi)$
Joint distribution between *z* and *x*

1) Monte Carlo approximation of ∇_{ϕ} ELBO(ϕ) via samples $z^1, \ldots, z^B \sim q_{\phi}$





Gradient-based variational inference:

ELBC Variational parameters ϕ

Will need gradients of the ELBO

What we will cover:

Monte Carlo approximations of the gradient

- 1) Monte Carlo approximation of ∇_{ϕ} ELBO(ϕ) via samples $z^1, \ldots, z^B \sim q_{\phi}$
- 2) Score function gradient estimator
- 3) Reparameterization gradient estimator

(uses ideas from VI + deep generative modeling)

Variational Bayes: learn an approximation for $\pi(z) = p(z | x)$ for data x

latent variables *z*.

$$D(\phi) = \int \log \left(\frac{p_{\theta}(z, x)}{q_{\phi}(z)} \right) q_{\phi}(z) dz$$

: $\nabla_{\phi} \text{ELBO}(\phi)$
Joint distribution between *z* and *x*

Amortized inference: learn an approximation for p(z | x) for general x







Mohamed et al.. Monte Carlo Gradient Estimation in Machine Learning. Journal of Machine Learning Research, 2020.

Example: ELBO(ϕ) = $\mathbb{E}_{q_{\phi}(z)} \left| \log \left(\frac{p_{\theta}(z, x)}{q_{\phi}(z)} \right) \right|$







Other applications: reinforcement learning, experimental design

Mohamed et al.. Monte Carlo Gradient Estimation in Machine Learning. Journal of Machine Learning Research, 2020.

Example: ELBO(ϕ) = $\mathbb{E}_{q_{\phi}(z)} \left| \log \left(\frac{p_{\theta}(z, x)}{q_{\phi}(z)} \right) \right|$









- **Example:** ELBO(ϕ) = $\mathbb{E}_{q_{\phi}(z)} \left| \log \left(\frac{p_{\theta}(z, x)}{q_{\phi}(z)} \right) \right|$
- Other applications: reinforcement learning, experimental design
- **Goal:** learn the distributional parameters ϕ via gradient descent. Need gradients: $\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)}[f(z;\phi,\theta)]$







- **Example:** ELBO(ϕ) = $\mathbb{E}_{q_{\phi}(z)} \left| \log \left(\frac{p_{\theta}(z, x)}{q_{\phi}(z)} \right) \right|$
- Other applications: reinforcement learning, experimental design
- **Goal:** learn the distributional parameters ϕ via gradient descent. Need gradients: $\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)}[f(z;\phi,\theta)]$ **But:** typically this expectation is intractable!







Two strategies that appear frequently in approximate inference: 1) Score function gradient estimator 2) Reparameterization gradient estimator

- **Example:** ELBO(ϕ) = $\mathbb{E}_{q_{\phi}(z)} \left| \log \left(\frac{p_{\theta}(z, x)}{q_{\phi}(z)} \right) \right|$
- Other applications: reinforcement learning, experimental design
- **Goal:** learn the distributional parameters ϕ via gradient descent. Need gradients: $\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)}[f(z;\phi,\theta)]$ **But:** typically this expectation is intractable!







Two strategies that appear frequently in approximate inference:

- 1) Score function gradient estimator
- 2) Reparameterization gradient estimator

- **Example:** ELBO(ϕ) = $\mathbb{E}_{q_{\phi}(z)} \left| \log \left(\frac{p_{\theta}(z, x)}{q_{\phi}(z)} \right) \right|$
- Other applications: reinforcement learning, experimental design
- **Goal:** learn the distributional parameters ϕ via gradient descent. Need gradients: $\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)}[f(z;\phi,\theta)]$ **But:** typically this expectation is intractable!

- We will also discuss two "black-box" strategies for VI based on these two estimators
 - Mohamed et al.. Monte Carlo Gradient Estimation in Machine Learning. Journal of Machine Learning Research, 2020.







Monte Carlo estimates

$\mathscr{L}(\phi) := \mathbb{E}_{q_{\phi}(z)}[f(z; \phi, \theta)]$

Example: ELBO(ϕ) = $\mathbb{E}_{q_{\phi}(z)} \left| \log \left(\frac{p_{\theta}(z, x)}{q_{\phi}(z)} \right) \right|$



Monte Carlo estimates

Monte Carlo estimate: sample $z^b \sim q_{\phi}(z)$

 $\mathscr{L}(\phi) := \mathbb{E}_{q_{\phi}(z)}[f(z;\phi,\theta)] \qquad \text{Example: ELBO}(\phi) = \mathbb{E}_{q_{\phi}(z)}\left|\log\left(\frac{p_{\theta}(z,x)}{q_{\phi}(z)}\right)\right|$



Monte Carlo estimates

$$\mathscr{L}(\phi) := \mathbb{E}_{q_{\phi}(z)}[f(z; \phi, \theta)]$$

Monte Carlo estimate: sample $z^b \sim q_{\phi}(z)$ $\mathscr{L}(\phi) \approx \frac{1}{R} \sum_{b=1}^{B} f(z^{b}; \phi, \theta)$

Example: ELBO(ϕ) = $\mathbb{E}_{q_{\phi}(z)} \left| \log \left(\frac{p_{\theta}(z, x)}{q_{\phi}(z)} \right) \right|$



Monte Carlo estimates **Example:** ELBO(ϕ) = $\mathbb{E}_{q_{\phi}(z)} \left| \log \left(\frac{p_{\theta}(z, x)}{q_{\phi}(z)} \right) \right|$

$$\mathscr{L}(\phi) := \mathbb{E}_{q_{\phi}(z)}[f(z;\phi,\theta)]$$

Monte Carlo estimate: sample $z^b \sim q_{\phi}(z)$ $\mathscr{L}(\phi) \approx \frac{1}{R} \sum_{b=1}^{B} f(z^{b}; \phi, \theta)$

How to estimate the gradient $\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)}[f(z; \phi, \theta)]$ using Monte Carlo?



Monte Carlo estimates **Example:** ELBO(ϕ) = $\mathbb{E}_{q_{\phi}(z)} \left| \log \left(\frac{p_{\theta}(z, x)}{q_{\phi}(z)} \right) \right|$ $q_{\phi}(z)$

$$\begin{aligned} \mathscr{L}(\phi) &:= \mathbb{E}_{q_{\phi}(z)}[f(z;\phi,\theta)] \\ \text{Monte Carlo estimate: sample } z^b \sim q \\ \mathscr{L}(\phi) &\approx \frac{1}{B} \sum_{b=1}^{B} f(z^b;\phi,\theta) \end{aligned}$$

How to estimate the gradient $\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)}[f(z; \phi, \theta)]$ using Monte Carlo?

Desirable properties of estimators:

- Unbiased: E[estimator] = true value
- Consistent: recover true value as $B \rightarrow \infty$
- we prefer the lower variance once
- Fast computation

• Low variance: if we have two estimators that use the same amount of computation,



a.k.a. likelihood ratio method, REINFORCE estimator, "log derivative trick"



a.k.a. likelihood ratio method, REINFORCE estimator, "log derivative trick"

Score function: (a.k.a. "Fisher" score)

 $\nabla_{\phi} \log q_{\phi}(z) = \frac{\nabla_{\phi} q_{\phi}(z)}{q_{\phi}(z)}$



a.k.a. likelihood ratio method, REINFORCE estimator, "log derivative trick"

Score function: (a.k.a. "Fisher" score)

 $\nabla_{\phi} \log q_{\phi}(z) = \frac{\nabla_{\phi} q_{\phi}(z)}{q_{\phi}(z)}$

(Note: NOT the "Stein" score function $\nabla_z \log q_{\phi}(z)$)



a.k.a. likelihood ratio method, REINFORCE estimator, "log derivative trick"

Score function: (a.k.a. "Fisher" score)

$$\nabla_{\phi} \log q_{\phi}(z) = \frac{\nabla_{\phi} q_{\phi}(z)}{q_{\phi}(z)}$$

Gradient of the expectation:

(Note: NOT the "Stein" score function $\nabla_z \log q_{\phi}(z)$)



a.k.a. likelihood ratio method, REINFORCE estimator, "log derivative trick"

Score function: (a.k.a. "Fisher" score)

$$\nabla_{\phi} \log q_{\phi}(z) = \frac{\nabla_{\phi} q_{\phi}(z)}{q_{\phi}(z)}$$

Gradient of the expectation:

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)}[f(z)] = \nabla_{z}$$

(Note: NOT the "Stein" score function $\nabla_z \log q_{\phi}(z)$)

 $7_{\phi} \left[q_{\phi}(z) f(z) dz \right]$


a.k.a. likelihood ratio method, REINFORCE estimator, "log derivative trick"

Score function: (a.k.a. "Fisher" score)

$$\nabla_{\phi} \log q_{\phi}(z) = \frac{\nabla_{\phi} q_{\phi}(z)}{q_{\phi}(z)}$$

Gradient of the expectation:

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)}[f(z)] = \nabla_{z}$$

(Note: NOT the "Stein" score function $\nabla_z \log q_{\phi}(z)$)

 $7_{\phi} q_{\phi}(z) f(z) dz$

Assume you can swap derivatives and integrals



a.k.a. likelihood ratio method, REINFORCE estimator, "log derivative trick"

Score function: (a.k.a. "Fisher" score)

$$\nabla_{\phi} \log q_{\phi}(z) = \frac{\nabla_{\phi} q_{\phi}(z)}{q_{\phi}(z)}$$

Gradient of the expectation:

 $\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)}[f(z)] = \nabla_{\phi} \int q_{\phi}(z) f(z) dz$ $= \int \nabla_{\phi} q_{\phi}(z) f(z) \qquad \text{Assume you can swap} \\ \text{derivatives and integrals}$

(Note: NOT the "Stein" score function $\nabla_z \log q_{\phi}(z)$)



a.k.a. likelihood ratio method, REINFORCE estimator, "log derivative trick"

Score function: (a.k.a. "Fisher" score)

$$\nabla_{\phi} \log q_{\phi}(z) = \frac{\nabla_{\phi} q_{\phi}(z)}{q_{\phi}(z)}$$

Gradient of the expectation:

 $\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)}[f(z)] = \nabla_{\phi} \int q_{\phi}(z) f(z) dz$ $= \int \nabla_{\phi} q_{\phi}(z) f(z) \qquad \text{Assume you can swap} \\ \text{derivatives and integrals}$

(Note: NOT the "Stein" score function $\nabla_z \log q_{\phi}(z)$)



a.k.a. likelihood ratio method, REINFORCE estimator, "log derivative trick"

Score function: (a.k.a. "Fisher" score)

$$\nabla_{\phi} \log q_{\phi}(z) = \frac{\nabla_{\phi} q_{\phi}(z)}{q_{\phi}(z)}$$

Gradient of the expectation:

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)}[f(z)] = \nabla_{a}$$

"Log derivative trick" $\nabla_{\phi} q_{\phi}(z) = q_{\phi}(z) \nabla_{\phi} \log q_{\phi}(z)$ (Note: NOT the "Stein" score function $\nabla_z \log q_{\phi}(z)$)

 $= \nabla_{\phi} \int q_{\phi}(z) f(z) dz$ = $\int \nabla_{\phi} q_{\phi}(z) f(z)$ Assume you can swap derivatives and integrals



a.k.a. likelihood ratio method, REINFORCE estimator, "log derivative trick"

Score function: (a.k.a. "Fisher" score)

$$\nabla_{\phi} \log q_{\phi}(z) = \frac{\nabla_{\phi} q_{\phi}(z)}{q_{\phi}(z)}$$

Gradient of the expectation:

 $\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)}[f(z)] = \nabla_{\phi} \int q_{\phi}(z) f(z) dz$

"Log derivative trick" $\nabla_{\phi} q_{\phi}(z) = q_{\phi}(z) \nabla_{\phi} \log q_{\phi}(z)$ (Note: NOT the "Stein" score function $\nabla_z \log q_{\phi}(z)$)

 $= \int \nabla_{\phi} q_{\phi}(z) f(z) \qquad \text{Assume you can swap} \\ \text{derivatives and integrals} \\ = \int q_{\phi}(z) f(z) \nabla_{\phi} \log q_{\phi}(z) dz$



a.k.a. likelihood ratio method, REINFORCE estimator, "log derivative trick"

Score function: (a.k.a. "Fisher" score)

$$\nabla_{\phi} \log q_{\phi}(z) = \frac{\nabla_{\phi} q_{\phi}(z)}{q_{\phi}(z)}$$

Gradient of the expectation:

 $\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)}[f(z)] = \nabla_{\phi} \left[q_{\phi}(z) f(z) dz \right]$

"Log derivative trick" $\nabla_{\phi} q_{\phi}(z) = q_{\phi}(z) \nabla_{\phi} \log q_{\phi}(z)$ (Note: NOT the "Stein" score function $\nabla_z \log q_{\phi}(z)$)

 $= \int \nabla_{\phi} q_{\phi}(z) f(z) \qquad \text{Assume you can swap} \\ \text{derivatives and integrals} \\ = \int q_{\phi}(z) f(z) \nabla_{\phi} \log q_{\phi}(z) dz$



a.k.a. likelihood ratio method, REINFORCE estimator, "log derivative trick"

Score function: (a.k.a. "Fisher" score)

$$\nabla_{\phi} \log q_{\phi}(z) = \frac{\nabla_{\phi} q_{\phi}(z)}{q_{\phi}(z)}$$

Gradient of the expectation:

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)}[f(z)] = \nabla_{\phi} \mathbb{E}_{q_{\phi}(z)}[f(z)] = \nabla_{\phi$$

"Log derivative trick" $\nabla_{\phi} q_{\phi}(z) = q_{\phi}(z) \nabla_{\phi} \log q_{\phi}(z)$

Sample
$$z^b \sim q_{\phi}(z)$$

(Note: NOT the "Stein" score function $\nabla_z \log q_{\phi}(z)$)

 $= \nabla_{\phi} \int q_{\phi}(z) f(z) dz$ = $\int \nabla_{\phi} q_{\phi}(z) f(z) f(z)$ Assume you can swap derivatives and integrals = $\int q_{\phi}(z) f(z) \nabla_{\phi} \log q_{\phi}(z) dz$



a.k.a. likelihood ratio method, REINFORCE estimator, "log derivative trick"

Score function: (a.k.a. "Fisher" score)

$$\nabla_{\phi} \log q_{\phi}(z) = \frac{\nabla_{\phi} q_{\phi}(z)}{q_{\phi}(z)}$$

Gradient of the expectation:

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)}[f(z)] = \nabla$$

"Log derivative trick" $\nabla_{\phi} q_{\phi}(z) = q_{\phi}(z) \nabla_{\phi} \log q_{\phi}(z)$

Sample
$$z^b \sim q_{\phi}(z)$$

(Note: NOT the "Stein" score function $\nabla_z \log q_{\phi}(z)$)

 $| = \nabla_{\phi} \int q_{\phi}(z) f(z) dz$ = $\int \nabla_{\phi} q_{\phi}(z) f(z) f(z)$ Assume you can swap derivatives and integrals

 $= \int q_{\phi}(z) f(z) \nabla_{\phi} \log q_{\phi}(z) dz$

 $\approx \frac{1}{R} \sum \nabla_{\phi} \log q_{\phi}(z^b) f(z^b)$ *b*=1



a.k.a. likelihood ratio method, REINFORCE estimator, "log derivative trick"

Score function: (a.k.a. "Fisher" score)

$$\nabla_{\phi} \log q_{\phi}(z) = \frac{\nabla_{\phi} q_{\phi}(z)}{q_{\phi}(z)}$$

Gradient of the expectation:

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)}[f(z)] = \nabla$$

"Log derivative trick" $\nabla_{\phi} q_{\phi}(z) = q_{\phi}(z) \nabla_{\phi} \log q_{\phi}(z)$

Sample
$$z^b \sim q_{\phi}(z)$$

(Note: NOT the "Stein" score function $\nabla_z \log q_{\phi}(z)$)

 $= \nabla_{\phi} \int q_{\phi}(z) f(z) dz$ = $\int \nabla_{\phi} q_{\phi}(z) f(z) \qquad \text{Assume you can swap} \\ \text{derivatives and integrals}$ = $\int q_{\phi}(z) f(z) \nabla_{\phi} \log q_{\phi}(z) dz$

 $\sum \nabla_{\phi} \log q_{\phi}(z^{b}) f(z^{b}) \quad \begin{array}{l} \text{Variance reduction} \\ \text{("control variate")} \end{array}$ $\approx \frac{1}{B} \sum_{n=1}^{\infty}$ ("control variate") b=1 $[f(z^b) - \beta]$



Approach 2: reparameterization gradient estimator

Sampling paths: suppose we can alternatively generate samples as follows:

a.k.a. pathwise gradient estimator

 $\hat{z} \sim q_{\phi}(z;\phi) \equiv \hat{z} = g(\hat{\epsilon};\phi), \quad \hat{\epsilon} \sim p(\epsilon)$



Approach 2: reparameterization gradient estimator

Sampling paths: suppose we can alternatively generate samples as follows:

 $\hat{z} \sim q_{\phi}(z;\phi) \equiv \hat{z}$

a.k.a. pathwise gradient estimator

$$\hat{\epsilon} = g(\hat{\epsilon}; \phi), \quad \hat{\epsilon} \sim p(\epsilon)$$

deterministic path



Approach 2: reparameterization gradient estimator

Sampling paths: suppose we can alternatively generate samples as follows:

 $\hat{z} \sim q_{\phi}(z;\phi) \equiv \hat{z}$

deterministic path

a.k.a. pathwise gradient estimator

$$\hat{\epsilon} = g(\hat{\epsilon}; \phi), \quad \hat{\epsilon} \sim p(\epsilon)$$

base distribution that is independent of ϕ



Sampling paths: suppose we can alternatively generate samples as follows:

$$\hat{z} \sim q_{\phi}(z;\phi) \equiv \hat{z} = g(\hat{\epsilon};\phi), \quad \hat{\epsilon} \sim p(\epsilon)$$

- **Example:** Multivariate Gaussian $p(z; \phi) = \mathcal{N}(z \mid \mu, \Sigma)$,
 - base is $p(\epsilon) = \mathcal{N}(0,I); g(\epsilon)$

deterministic path base distribution that is independent of ϕ

$$f(\epsilon; \theta) = \mu + L\epsilon$$
, where $\Sigma = LL^{\mathsf{T}}$



Sampling paths: suppose we can alternatively generate samples as follows:

$$\hat{z} \sim q_{\phi}(z;\phi) \equiv \hat{z} = g(\hat{\epsilon};\phi), \quad \hat{\epsilon} \sim p(\epsilon)$$

- **Example:** Multivariate Gaussian $p(z; \phi) = \mathcal{N}(z \mid \mu, \Sigma)$,

deterministic path base distribution that is independent of ϕ

base is $p(\epsilon) = \mathcal{N}(0,I); \quad g(\epsilon;\theta) = \mu + L\epsilon$, where $\Sigma = LL^{\top}$

"Reparameterization trick": $\mathbb{E}_{q_{\phi}(z)}[f(z)] = \mathbb{E}_{p(\epsilon)}[f(g(\epsilon; \phi))]$



Sampling paths: suppose we can alternatively generate samples as follows:

$$\hat{z} \sim q_{\phi}(z;\phi) \equiv \hat{z} = g(\hat{\epsilon};\phi), \quad \hat{\epsilon} \sim p(\epsilon)$$

base distribution that is independent of ϕ deterministic path **Example:** Multivariate Gaussian $p(z; \phi) = \mathcal{N}(z \mid \mu, \Sigma)$,

base is $p(\epsilon) = \mathcal{N}(0,I); \quad g(\epsilon;\theta) = \mu + L\epsilon$, where $\Sigma = LL^+$

Application to gradient estimator:

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)}[f(z;$$

"Reparameterization trick": $\mathbb{E}_{q_{\phi}(z)}[f(z)] = \mathbb{E}_{p(\epsilon)}[f(g(\epsilon; \phi))]$



Sampling paths: suppose we can alternatively generate samples as follows:

$$\hat{z} \sim q_{\phi}(z;\phi) \equiv \hat{z} = g(\hat{\epsilon};\phi), \quad \hat{\epsilon} \sim p(\epsilon)$$

base distribution that is independent of ϕ deterministic path **Example:** Multivariate Gaussian $p(z; \phi) = \mathcal{N}(z \mid \mu, \Sigma)$,

base is $p(\epsilon) = \mathcal{N}(0,I); g$

"Reparameterization trick":

Application to gradient estimator:

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)}[f(z;$$

$$(\epsilon; \theta) = \mu + L\epsilon$$
, where $\Sigma = LL^{\mathsf{T}}$

$$\mathbb{E}_{q_{\phi}(z)}[f(z)] = \mathbb{E}_{p(\epsilon)}[f(g(\epsilon; \phi))]$$

$\phi)] = \mathbb{E}_{p(\epsilon)} [\nabla_{\phi} f(g(\epsilon; \phi))]$



Sampling paths: suppose we can alternatively generate samples as follows:

$$\hat{z} \sim q_{\phi}(z;\phi) \equiv \hat{z} = g(\hat{\epsilon};\phi), \quad \hat{\epsilon} \sim p(\epsilon)$$

deterministic path base distribution that is independent of ϕ **Example:** Multivariate Gaussian $p(z; \phi) = \mathcal{N}(z \mid \mu, \Sigma)$,

base is
$$p(\epsilon) = \mathcal{N}(0,I)$$
; $g(\epsilon;\theta) = \mu + L\epsilon$, where $\Sigma = LL^{\top}$

"Reparameterization trick":

Application to gradient estimator:

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)}[f(z;$$

Draw samples $\hat{\epsilon}^b \sim p(\epsilon)$

$$\mathbb{E}_{q_{\phi}(z)}[f(z)] = \mathbb{E}_{p(\epsilon)}[f(g(\epsilon; \phi))]$$

$\phi)] = \mathbb{E}_{p(\epsilon)}[\nabla_{\phi}f(g(\epsilon;\phi))]$



Sampling paths: suppose we can alternatively generate samples as follows:

$$\hat{z} \sim q_{\phi}(z;\phi) \equiv \hat{z} = g(\hat{\epsilon};\phi), \quad \hat{\epsilon} \sim p(\epsilon)$$

base distribution that is independent of ϕ deterministic path **Example:** Multivariate Gaussian $p(z; \phi) = \mathcal{N}(z \mid \mu, \Sigma)$,

base is
$$p(\epsilon) = \mathcal{N}(0,I)$$
; $g(\epsilon;\theta) = \mu + L\epsilon$, where $\Sigma = LL^{\top}$

"Reparameterization trick": \mathbb{E}

Application to gradient estimator:

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)}[f(z;\phi)] = \mathbb{E}_{p(\epsilon)}[\nabla_{\phi} f(g(\epsilon;\phi))]$$
$$\approx \frac{1}{B} \sum_{b=1}^{B} \nabla_{\phi} f(g(\hat{\epsilon}^{b};\phi))$$

Draw samples $\hat{\epsilon}^b \sim p(\epsilon)$

$$\mathbb{E}_{q_{\phi}(z)}[f(z)] = \mathbb{E}_{p(\epsilon)}[f(g(\epsilon; \phi))]$$



Comparison of estimators $q(z; \phi) = \mathcal{N}(z \,|\, \mu, \sigma^2)$ $f(z) = (z - k)^2$ $\phi \in \{\mu, \sigma\}$



Mohamed et al. Monte Carlo Gradient Estimation in Machine Learning. Journal of Machine Learning Research, 2020.





Comparison of estimators $q(z; \phi) = \mathcal{N}(z \,|\, \mu, \sigma^2)$ $f(z) = (z - k)^2$ $\phi \in \{\mu, \sigma\}$



Mohamed et al. Monte Carlo Gradient Estimation in Machine Learning. Journal of Machine Learning Research, 2020.





Comparison of estimators $q(z; \phi) = \mathcal{N}(z \,|\, \mu, \sigma^2)$ $\phi \in \{\mu, \sigma\}$



Mohamed et al. Monte Carlo Gradient Estimation in Machine Learning. Journal of Machine Learning Research, 2020.



Score function Score + variance reduction





Comparison of estimators $q(z; \phi) = \mathcal{N}(z \,|\, \mu, \sigma^2)$ $f(z) = (z - k)^2$ $\phi \in \{\mu, \sigma\}$



Mohamed et al. Monte Carlo Gradient Estimation in Machine Learning. Journal of Machine Learning Research, 2020.





Comparison of estimators $q(z; \phi) = \mathcal{N}(z | \mu, \sigma^2)$ $\phi \in \{\mu, \sigma\}$ $f(z) = (z - k)^2$



Score function: more general; more sensitive to dimensionality; generally higher variance **Reparameterization:** needs a differentiable cost function; variance more well behaved

Mohamed et al. Monte Carlo Gradient Estimation in Machine Learning. Journal of Machine Learning Research, 2020.





Uses the score function gradient estimator of the ELBO Sample $z^1, ..., z^B \sim q_{\phi}(z)$ Set the step size η Gradient update $\phi \leftarrow \phi + \eta \frac{1}{B} \sum_{b=1}^{B} \nabla_{\phi} \log q_{\phi}(z^{b}) \log \left(\frac{p(z^{b}, x)}{q_{\phi}(z^{b})}\right)$ b=1

Ranganath, Gerrish, Blei. Black box variational inference. AISTATS, 2014.



Uses the score function gradient estimator of the ELBO Sample $z^1, ..., z^B \sim q_{\phi}(z)$ Set the step size η Gradient update $\phi \leftarrow \phi + \eta \frac{1}{B} \sum_{\phi} \nabla_{\phi} \log$ b=1

Ranganath, Gerrish, Blei. Black box variational inference. AISTATS, 2014.

$$g q_{\phi}(z^b) \log\left(\frac{p(z^b, x)}{q_{\phi}(z^b)}\right)$$



Uses the score function gradient estimator of the ELBO Sample $z^1, ..., z^B \sim q_{\phi}(z)$ Set the step size η Gradient update $\phi \leftarrow \phi + \eta \frac{1}{B} \sum_{\phi}^{B} \nabla_{\phi} \log$ b=1

Ranganath, Gerrish, Blei. Black box variational inference. AISTATS, 2014.

$$g q_{\phi}(z^b) \log\left(\frac{p(z^b, x)}{q_{\phi}(z^b)}\right)$$

"Black box": few assumptions are needed on p(z, x)

- Sample from q
- Evaluate score
- Evaluate joint





Uses the score function gradient estimator of the ELBO Sample $z^1, ..., z^B \sim q_{\phi}(z)$ Set the step size η Gradient update $\phi \leftarrow \phi + \eta \frac{1}{B} \sum_{\phi} \nabla_{\phi} \log$ b=1

Black-box VI is used with techniques to reduce the variance of the estimator (Rao-Blackwellization and control variates)

Ranganath, Gerrish, Blei. Black box variational inference. AISTATS, 2014.

$$gq_{\phi}(z^b)\log\left(\frac{p(z^b,x)}{q_{\phi}(z^b)}\right)$$

"Black box": few assumptions are needed on p(z, x)

- Sample from q
- Evaluate score
- Evaluate joint





Uses the score function gradient estimator of the ELBO Sample $z^1, ..., z^B \sim q_{\phi}(z)$ Set the step size η Gradient update $\phi \leftarrow \phi + \eta \frac{1}{B} \sum_{\phi}^{B} \nabla_{\phi} \log$ h=1

Black-box VI is used with techniques to reduce the variance of the estimator (Rao-Blackwellization and control variates)

Ranganath, Gerrish, Blei. Black box variational inference. AISTATS, 2014.

$$g q_{\phi}(z^b) \log\left(\frac{p(z^b, x)}{q_{\phi}(z^b)}\right)$$

"Black box": few assumptions are needed on p(z, x)

- Sample from q
- Evaluate score
- Evaluate joint







Transform target π to one on real coordinate space $\tilde{\pi}$ (via change of variables formula)



(a) Latent variable space

Kucukelbir, et al. Automatic differentiation variational inference. Journal of Machine Learning Research, 2018.





Transform target π to one on real coordinate space $\tilde{\pi}$ (via change of variables formula)



(a) Latent variable space (b) Real coordinate space Uses the reparameterization gradient estimator of the ELBO:

Kucukelbir, et al. Automatic differentiation variational inference. Journal of Machine Learning Research, 2018.





Transform target π to one on real coordinate space $\tilde{\pi}$ (via change of variables formula)



(a) Latent variable space

Uses the reparameterization gradient estimator of the ELBO: Variational family is a Gaussian with μ and Σ (diagonal or full-covariance)

Kucukelbir, et al. Automatic differentiation variational inference. Journal of Machine Learning Research, 2018.





Transform target π to one on real coordinate space $\tilde{\pi}$ (via change of variables formula)



(a) Latent variable space

Uses the reparameterization gradient estimator of the ELBO: Variational family is a Gaussian with μ and Σ (diagonal or full-covariance) Base distribution is $p(\epsilon) = \mathcal{N}(0,I)$ and e.g., $g(\epsilon; \theta) = \mu + L\epsilon$, where $LL^{\perp} = \Sigma$

Kucukelbir, et al. Automatic differentiation variational inference. Journal of Machine Learning Research, 2018.





Transform target π to one on real coordinate space $\tilde{\pi}$ (via change of variables formula)



(a) Latent variable space

Uses the reparameterization gradient estimator of the ELBO: Variational family is a Gaussian with μ and Σ (diagonal or full-covariance) Base distribution is $p(\epsilon) = \mathcal{N}(0,I)$ and e.g., $g(\epsilon; \theta) = \mu + L\epsilon$, where $LL^{\perp} = \Sigma$ Requires that the log joint distribution is differentiable (i.e., score $V_{\phi} \log p(z, x)$).

Kucukelbir, et al. Automatic differentiation variational inference. Journal of Machine Learning Research, 2018.





Comparison of BBVI vs ADVI variance

score function gradient vs reparameterization gradient



Kucukelbir, Tran, Ranganath, Gelman, Blei. Automatic differentiation variational inference. Journal of Machine Learning Research, 2018.





Summary: score function vs reparameterization gradients

Score function gradient estimator

- More general: doesn't need differentiable model $\rho(z) = p(z, x)$ Applies to continuous and discrete distributions • High variance: needs to be used with variance reduction techniques

Reparameterization gradients

- Requires differentiable models
- Needs to be "reparameterizable": i.e
- Generally controlled variance (less s

s.,
$$z = g(\hat{\epsilon}, \phi)$$

sensitive to dimension)



Software implementations & other advances

Black-box VI and ADVI: Implemented in many software systems


Software implementations & other advances Black-box VI and ADVI: Implemented in many software systems Probabilistic programming software: e.g., Stan, Pyro, Turing, PyMC





Software implementations & other advances Black-box VI and ADVI: Implemented in many software systems Probabilistic programming software: e.g., Stan, Pyro, Turing, PyMC 5 **Deterministic ADVI (DADVI):** swaps differentiation and sampling; lower-variance, improved stability; alternative solvers

Giordano, et al. Black Box Variational Inference with a Deterministic Objective: Faster, More Accurate, and Even More Black Box. JMLR 2024.







Software implementations & other advances

Black-box VI and ADVI: Implemented in many software systems

Probabilistic programming software: e.g., Stan, Pyro, Turing, PyMC

improved stability; alternative solvers

Implicit reparameterization: based on implicit differentiation;

- allows reparameterization to be applied to a broader class of distributions, e.g., beta, gamma, and Dirichlet
- implemented in, e.g., PyTorch, JAX (numpyro, distrax)

Giordano, et al. Black Box Variational Inference with a Deterministic Objective: Faster, More Accurate, and Even More Black Box. JMLR 2024. Figurnov et al. Implicit reparameterization gradients. *NeurIPS* 2018.



- **Deterministic ADVI (DADVI):** swaps differentiation and sampling; lower-variance,





Software implementations & other advances

Black-box VI and ADVI: Implemented in many software systems

Probabilistic programming software: e.g., Stan, Pyro, Turing, PyMC 5

improved stability; alternative solvers

Implicit reparameterization: based on implicit differentiation;

- allows reparameterization to be applied to a broader class of distributions, e.g., beta, gamma, and Dirichlet
- implemented in, e.g., PyTorch, JAX (numpyro, distrax)

Scaling to large datasets: can subsample data to speed up VI, e.g.,

$$\sum_{n=1}^{N} \log p(x_n | z) \approx \frac{N}{M} \sum_{m \in \mathcal{M}} \log p(x_m | z)$$

Giordano, et al. Black Box Variational Inference with a Deterministic Objective: Faster, More Accurate, and Even More Black Box. JMLR 2024. Figurnov et al. Implicit reparameterization gradients. *NeurIPS* 2018.



- **Deterministic ADVI (DADVI):** swaps differentiation and sampling; lower-variance,





Software implementations & other advances

Black-box VI and ADVI: Implemented in many software systems

Probabilistic programming software: e.g., Stan, Pyro, Turing, PyMC 5

improved stability; alternative solvers

Implicit reparameterization: based on implicit differentiation;

- allows reparameterization to be applied to a broader class of distributions, e.g., beta, gamma, and Dirichlet
- implemented in, e.g., PyTorch, JAX (numpyro, distrax)

Scaling to large datasets: can subsample data to speed up VI, e.g.,

$$\sum_{n=1}^{N} \log p(x_n | z) \approx \frac{N}{M} \sum_{m \in M} \sum_{m \in M}$$

Giordano, et al. Black Box Variational Inference with a Deterministic Objective: Faster, More Accurate, and Even More Black Box. JMLR 2024. Figurnov et al. Implicit reparameterization gradients. *NeurIPS* 2018.



- **Deterministic ADVI (DADVI):** swaps differentiation and sampling; lower-variance,

A different approach: $\sum \log p(x_m | z)$ amortization M







Bayesian inference

Fitting a probabilistic model p(z | x) for specific input x

- Typically relies on knowing the likelihood

• Need to re-do inference every time we see a new data set x (expensive)





Bayesian inference Fitting a probabilistic model p(z | x) for specific input x

- Typically relies on knowing the likelihood

Amortized inference: draws from training/test ideas in ML but for the posterior

• Need to re-do inference every time we see a new data set x (expensive)





Bayesian inference

- Need to re-do inference every time we see a new data set x (expensive) • Typically relies on knowing the likelihood

Amortized inference: draws from training/test ideas in ML but for the posterior







Fitting a probabilistic model p(z | x) for specific input x





Bayesian inference

- Need to re-do inference every time we see a new data set x (expensive) Typically relies on knowing the likelihood

Amortized inference: draws from training/test ideas in ML but for the posterior



Data *x*



Fitting a probabilistic model p(z | x) for specific input x

for *general* input x (slow)





Bayesian inference

- Need to re-do inference every time we see a new data set x (expensive) Typically relies on knowing the likelihood

Amortized inference: draws from training/test ideas in ML but for the posterior



Data *x*



Fitting a probabilistic model p(z | x) for specific input x





Bayesian inference

- Need to re-do inference every time we see a new data set x (expensive) Typically relies on knowing the likelihood

Amortized inference: draws from training/test ideas in ML but for the posterior

Data *x*





Applications:

- Simulation-based inference (e.g., neural posterior estimation)
- Generative modeling (e.g., variational autoencoders, normalizing flows) • Bayesian optimization (e.g., hyperparameter inference) • Meta-learning (e.g., neural processes)

Fitting a probabilistic model p(z | x) for specific input x





Example 1: variational autoencoders (VAEs)

Datapoint



Key ideas in VAEs:

- 1. Use the ELBO to approximate (lower bound) the log likelihood function
- 2. Amortized inference is used to approximate the posterior
- 3. Reparameterization trick used for training

Kingma and Welling. Auto-encoding variational Bayes. *ICLR* 2014. Rezende, Mohamed, Wierstra. Stochastic back propagation and approximate inference in deep generative models. *ICML* 2014.

ver bound) the log likelihood function proximate the posterior craining





pximation





pximation





pximation

"Decoder"













Instead, use approximate inference: $q_{\phi}(z | x) \approx p_{\theta}(z | x)$





Instead, use approximate inference: $q_{\phi}(z | x) \approx p_{\theta}(z | x)$

In VAEs, common to use a Gaussian, e.g.: $q_{\phi}(z \mid x) = \mathcal{N}(\Omega_{\mu}(x; \phi); \Omega_{\Sigma}(x; \phi))$





Instead, use approximate inference:

In VAEs, common to use a Gaussian, e.g.: $q_{\phi}(z \,|\, x) = \mathcal{N}(\Omega_{\mu}(x; \phi); \Omega_{\Sigma}(x; \phi))$

$$q_{\phi}(z \,|\, x) \approx p_{\theta}(z \,|\, x)$$

"Encoder" or "Inference network"





Instead, use approximate inference:

In VAEs, common to use a Gaussian, e.g.: "Encoder" or "Inference network" $q_{\phi}(z \mid x) = \mathcal{N}(\Omega_{\mu}(x; \phi); \Omega_{\Sigma}(x; \phi))$ **Classical VI:** local latent variables for each data point x_{μ}

$$q_{\phi}(z \,|\, x) \approx p_{\theta}(z \,|\, x)$$





- Maximum likelihood of $p_{\theta}(x) = \int p_{\theta}(x, z) dz$ is intractable!
- Instead, use approximate inference:

In VAEs, common to use a Gaussian, e.g.: "Encoder" or "Inference network" $q_{\phi}(z \mid x) = \mathcal{N}(\Omega_{\mu}(x; \phi); \Omega_{\Sigma}(x; \phi))$ **Classical VI:** local latent variables for each data point x_n Amortization: shared parameter ϕ across all data points

$$q_{\phi}(z \,|\, x) \approx p_{\theta}(z \,|\, x)$$













 $\leq \log p_{\theta}(x)$





- maximize the marginal likelihood of $p_{\theta}(x)$
- Thus, maximizing ELBO w.r.t. θ, ϕ will simultaneously: - minimize the KL divergence between q_{ϕ} and p_{θ}



 $\mathsf{ELBO}(\theta, \phi) = \mathbb{E}_{q_{\phi}(z|x)} \left[\log \left[\frac{p_{\theta}(x, z)}{q_{\phi}(z|x)} \right] \right]$



 $\mathsf{ELBO}(\theta, \phi) = \mathbb{E}_{q_{\phi}(z|x)} \left[\log \left[\frac{p_{\theta}(x, z)}{q_{\phi}(z|x)} \right] \right]$





$$\begin{split} \mathsf{ELBO}(\theta,\phi) &= \mathbb{E}_{q_{\phi}(z|x)} \bigg[\log \bigg[\frac{p_{\theta}(x,z)}{q_{\phi}(z|x)} \bigg] \bigg] & \qquad \mathsf{Recall for a Gaussian } \mathcal{N}(\mu,\Sigma) \\ &z &= \mu + L\epsilon, \quad \Sigma = LL^{\mathsf{T}} \\ &p(\epsilon) &= \mathcal{N}(0,I) \\ &= \mathbb{E}_{p(\epsilon)} \bigg[\log \bigg[\frac{p_{\theta}(x,z)}{q_{\phi}(z|x)} \bigg] \bigg], \quad z = g(\epsilon,\phi,x) \end{split}$$







$$\mathsf{ELBO}(\theta, \phi) = \mathbb{E}_{q_{\phi}(z|x)} \left[\log \left[\frac{1}{q_{\phi}(z|x)} \right] \right]$$
$$= \mathbb{E}_{p(\epsilon)} \left[\log \left[\frac{p_{\phi}(z|x)}{q_{\phi}(z|x)} \right] \right]$$

Form Monte Carlo estimator of the ELBO $\widehat{\mathscr{L}}_{\theta,\phi}(x;\epsilon)$: $\hat{\epsilon} \sim p(\epsilon)$ $z = g(\hat{\epsilon}, \phi, x)$ $\widehat{\mathcal{L}}_{\theta, \phi}(x; \hat{\epsilon}) = \log p_{\theta}(x, z) - \log q_{\phi}(z \mid x)$









$$\mathsf{ELBO}(\theta, \phi) = \mathbb{E}_{q_{\phi}(z|x)} \left[\log \left[\frac{p_{\phi}(z|x)}{q_{\phi}(z|x)} \right] \right]$$
$$= \mathbb{E}_{p(\epsilon)} \left[\log \left[\frac{p_{\phi}(z|x)}{q_{\phi}(z|x)} \right] \right]$$

Form Monte Carlo estimator of the ELBO $\widehat{\mathscr{L}}_{\theta,\phi}(x;\epsilon)$: $\hat{\epsilon} \sim p(\epsilon)$ $z = g(\hat{\epsilon}, \phi, x)$ $\widehat{\mathscr{L}}_{\theta, \phi}(x; \hat{\epsilon}) = \log p_{\theta}(x, z) - \log q_{\phi}(z \mid x)$





 $\nabla_{\theta,\phi} \mathsf{ELBO}(\theta,\phi) = \mathbb{E}_{p(\epsilon)} [\nabla_{\theta,\phi} \widehat{\mathscr{L}}_{\theta,\phi}(x;\epsilon)]$





Applications of VAEs

Generating MNIST digits:

6	¢.,	9		(, f)	7	Ŋ	¢,	0	3
fert.	Ņ	4	ст.	01	4	àr.	N. A.	Qr ²	Ò
4	1		ŝ	Û,	4	*	100	4	Altern
Ċ.	9	تۇرىيى	7	9	Ŋ	4	tan - √	14	Ú.
6	ù,	Ŧ	0	0	-0	. €.)		()	e)xa
Q	Q.		6	1	9		7	7	4
0 -	6	2	9	7	ţ.	4	4	\circ	Ą
Ĵ,	Э	*****	1	4	1	Ç		4	0
and the second	2	\$	7	6	4	Ċ,	5	З	7
	-	4.200	ŝ	1	4	0	9	Q	8



Rezende et al. Stochastic back propagation and approximate inference in deep generative models. *ICML* 2014. Kingma and Welling. Introduction to variational autoencoders. *Foundations and Trends® in Machine Learning*, 2019.

Image resynthesis:

Chemical design:





Example 2: Neural posterior estimation (NPE)

Goal: to compute the posterior with a test observation x_{obs} , i.e., $p(z | x_{obs})$

Papamakarios and Murray. Fast ε-free inference of simulation models with Bayesian conditional density estimation, 2018. Greenberg, Nonnenmacher, and Macke. Automatic posterior transformation for likelihood-free inference, 2019



Example 2: Neural posterior estimation (NPE)

Papamakarios and Murray. Fast ϵ -free inference of simulation models with Bayesian conditional density estimation, 2018. Greenberg, Nonnenmacher, and Macke. Automatic posterior transformation for likelihood-free inference, 2019

- **Goal:** to compute the posterior with a test observation x_{obs} , i.e., $p(z | x_{obs})$
- **But...** we do not have access to the likelihood. We can't evaluate $p(x \mid z)$



Example 2: Neural posterior estimation (NPE)

Papamakarios and Murray. Fast ϵ -free inference of simulation models with Bayesian conditional density estimation, 2018. Greenberg, Nonnenmacher, and Macke. Automatic posterior transformation for likelihood-free inference, 2019

- **Goal:** to compute the posterior with a test observation x_{obs} , i.e., $p(z | x_{obs})$
- **But...** we do not have access to the likelihood. We can't evaluate $p(x \mid z)$
 - **Examples:** expensive (differential equation); intractable (e.g. integral)


Suppose we can generate samples from $p(x \mid z)$

Papamakarios and Murray. Fast ϵ -free inference of simulation models with Bayesian conditional density estimation, 2018. Greenberg, Nonnenmacher, and Macke. Automatic posterior transformation for likelihood-free inference, 2019

- **Goal:** to compute the posterior with a test observation x_{obs} , i.e., $p(z | x_{obs})$
- **But...** we do not have access to the likelihood. We can't evaluate $p(x \mid z)$
 - **Examples:** expensive (differential equation); intractable (e.g. integral)



Suppose we can generate samples from $p(x \mid z)$

In NPE, perform simulation-based inference: $q_{\phi}(z | x_{obs}) \approx p(z | x_{obs})$

Papamakarios and Murray. Fast ϵ -free inference of simulation models with Bayesian conditional density estimation, 2018. Greenberg, Nonnenmacher, and Macke. Automatic posterior transformation for likelihood-free inference, 2019

- **Goal:** to compute the posterior with a test observation x_{obs} , i.e., $p(z | x_{obs})$
- **But...** we do not have access to the likelihood. We can't evaluate $p(x \mid z)$
 - **Examples:** expensive (differential equation); intractable (e.g. integral)



Suppose we can generate samples from $p(x \mid z)$

Papamakarios and Murray. Fast ϵ -free inference of simulation models with Bayesian conditional density estimation, 2018. Greenberg, Nonnenmacher, and Macke. Automatic posterior transformation for likelihood-free inference, 2019

- **Goal:** to compute the posterior with a test observation x_{obs} , i.e., $p(z | x_{obs})$
- **But...** we do not have access to the likelihood. We can't evaluate $p(x \mid z)$
 - **Examples:** expensive (differential equation); intractable (e.g. integral)
- In NPE, perform simulation-based inference: $q_{\phi}(z | x_{obs}) \approx p(z | x_{obs})$
 - Use *amortization* so that we can learn a $q_{\phi}(z \mid x)$ for general x



Suppose we can generate samples from $p(x \mid z)$

In NPE, perform simulation-based inference: $q_{\phi}(z | x_{obs}) \approx p(z | x_{obs})$

Key ideas in NPE:

- Minimize the expected forward KL via Monte Carlo
- 1. Amortized inference is used to approximate the posterior 2. Simulate pairs of the parameter and data 3.

Greenberg, Nonnenmacher, and Macke. Automatic posterior transformation for likelihood-free inference, 2019

- **Goal:** to compute the posterior with a test observation x_{obs} , i.e., $p(z | x_{obs})$
- **But...** we do not have access to the likelihood. We can't evaluate $p(x \mid z)$
 - **Examples:** expensive (differential equation); intractable (e.g. integral)

 - Use *amortization* so that we can learn a $q_{\phi}(z \mid x)$ for general x

Papamakarios and Murray. Fast ϵ -free inference of simulation models with Bayesian conditional density estimation, 2018.



NPE setup



NPE setup



NPE setup

Prior Simulator - can only generate samples (no likelihood)



Chose a flexible model $q_{\phi}(z \mid x)$, where ϕ represents parameters of a neural network

NPE setup

Prior Simulator - can only generate samples (no likelihood)





$$q_{\phi}(z \mid x) = \sum_{k=1}^{K} \Omega_{w_k}(z \mid x)$$

NPE setup

Prior Simulator - can only generate samples (no likelihood)

- Chose a flexible model $q_{\phi}(z \mid x)$, where ϕ represents parameters of a neural network Examples: Gaussian, mixture of Gaussians, normalizing flow
 - $(x;\phi) \mathcal{N}(z \mid \Omega_{\mu_{k}}(x;\phi), \Omega_{\Sigma_{k}}(x;\phi))$





$$q_{\phi}(z \mid x) = \sum_{k=1}^{K} \Omega_{w_{k}}(x; \phi) \mathcal{N}(z \mid \Omega_{\mu_{k}}(x; \phi), \Omega_{\Sigma_{k}}(x; \phi))$$

NPE setup

- Prior Simulator can only generate samples (no likelihood)
- Chose a flexible model $q_{\phi}(z \mid x)$, where ϕ represents parameters of a neural network Examples: Gaussian, mixture of Gaussians, normalizing flow

Want to minimize the *expected* forward KL divergence over data sets: $\mathscr{L}(\phi) = \int \mathscr{D}_{KL}(p(z \mid x); q_{\phi}(z; x)) p(x) dx$





$$q_{\phi}(z \mid x) = \sum_{k=1}^{K} \Omega_{w_{k}}(x; \phi) \mathcal{N}(z \mid \Omega_{\mu_{k}}(x; \phi), \Omega_{\Sigma_{k}}(x; \phi))$$

NPE setup

- Prior Simulator can only generate samples (no likelihood)
- Chose a flexible model $q_{\phi}(z \mid x)$, where ϕ represents parameters of a neural network Examples: Gaussian, mixture of Gaussians, normalizing flow

- Want to minimize the *expected* forward KL divergence over data sets: $\mathscr{L}(\phi) = \left[\mathscr{D}_{KL}(p(z \mid x); q_{\phi}(z; x)) p(x) \, dx \right]$
- **Question:** Why forward KL and not reverse KL? Problem: We can't evaluate p(z, x)







$$\mathscr{L}(\varphi) = \int \mathscr{D}_{KL}(p(z \mid x);$$

- **Goal:** to compute the posterior with a test observation x_{obs} , i.e., $p(z | x_{obs})$
 - $q_{\varphi}(z;x))p(x)dx$





$$\mathscr{L}(\varphi) = \int \mathscr{D}_{KL}(p(z|x); = \int \left[\int \log\left(\frac{p(z)}{q_{\varphi}(z)}\right) \right] dz$$

Goal: to compute the posterior with a test observation x_{obs} , i.e., $p(z | x_{obs})$

 $\frac{(z \mid x)}{(z;x)} \int p(z \mid x) dz \left| p(x) dx \right|$





$$\begin{aligned} \mathscr{L}(\varphi) &= \int \mathscr{D}_{KL}(p(z \mid x); \\ &= \int \left[\int \log \left(\frac{p(z)}{q_{\varphi}(z)} \right) \right] \\ &= \int \int \log \left(\frac{p(z \mid z)}{q_{\varphi}(z)} \right) \\ \end{aligned}$$

Goal: to compute the posterior with a test observation x_{obs} , i.e., $p(z | x_{obs})$

 $\frac{(z \mid x)}{(z \mid x)} \int p(z \mid x) dz \mid p(x) dx$ $\left(\frac{p(z \mid x)}{q_{\varphi}(z; x)}\right) p(z, x) dz dx$





$$\begin{aligned} \mathscr{L}(\varphi) &= \int \mathscr{D}_{KL}(p(z \mid x); \\ &= \int \left[\int \log \left(\frac{p(z)}{q_{\varphi}(z)} \right) \right] \\ &= \int \int \log \left(\frac{p(z \mid z)}{q_{\varphi}(z)} \right) \\ \end{aligned}$$

Goal: to compute the posterior with a test observation x_{obs} , i.e., $p(z | x_{obs})$

 $\frac{(z \mid x)}{(z \mid x)} \int p(z \mid x) dz \mid p(x) dx$ $\left(\frac{p(z \mid x)}{q_{\varphi}(z; x)}\right) p(z, x) dz dx$





$$\begin{aligned} \mathscr{L}(\varphi) &= \int \mathscr{D}_{KL}(p(z \mid x); \\ &= \int \left[\int \log \left(\frac{p(z)}{q_{\varphi}(z)} \right) \right] \\ \end{aligned}$$
The mate with arlo samples!

Approxim Monte Ca

Goal: to compute the posterior with a test observation x_{obs} , i.e., $p(z | x_{obs})$

 $\frac{z(x)}{z(x)} \int p(z(x)) dz p(x) dx$ $\frac{z(x)}{z(x)} \int p(z,x) dz dx$





$$\begin{aligned} \mathscr{L}(\varphi) &= \int \mathscr{D}_{KL}(p(z \mid x); \\ &= \int \left[\int \log \left(\frac{p(z \mid x)}{q_{\varphi}(z)} \right) \right] \\ \end{aligned}$$

$$\begin{aligned} &= \int \int \log \left(\frac{p(z \mid x)}{q_{\varphi}(z)} \right) \\ \end{aligned}$$

Approximate wit Monte Carlo sar

Draw
$$z^b \sim p(z)$$

 $x^b \sim p(x | z)$

Goal: to compute the posterior with a test observation x_{obs} , i.e., $p(z | x_{obs})$

 $\frac{z(x)}{z(x)} \int p(z(x)) dz \int p(x) dx$ $\left(\begin{array}{c} x \\ x \end{array} \right) \left(\begin{array}{c} p(z,x) \\ p(z,x) \\ z \\ x \end{array} \right) dz dx$





$$\begin{aligned} \mathscr{L}(\varphi) &= \int \mathscr{D}_{KL}(p(z \mid x); \\ &= \int \left[\int \log\left(\frac{p(z \mid x)}{q_{\varphi}(z)} \right) \right] \\ &= \int \int \log\left(\frac{p(z \mid x)}{q_{\varphi}(z)} \right) \\ &=$$

Approximate wi Monte Carlo sa

Draw $z^b \sim pb$ $x^b \sim p(x \,|\, z)$

Goal: to compute the posterior with a test observation x_{obs} , i.e., $p(z | x_{obs})$

 $\frac{z(x)}{z(x)} \int p(z(x)) dz \int p(x) dx$ $\frac{|x|}{|x|}$ p(z,x) dz dx

 $\hat{\mathscr{L}}(\varphi) = -\frac{1}{B}\sum \log q_{\varphi}(z^b; x^b)$ b=1





Input: prior and simulator p(z, x), conditional density estimator q_{ω} For b = 1, ..., BDraw $z^b, x^b \sim p(z, x)$

Train model to minimize the expected forward KL on the parameter-data pairs: $-\frac{1}{R}\sum \log q_{\varphi}(z^b;x^b)$ b=1

$$\varphi^* = \operatorname*{arg\,min}_{\varphi}$$

Output: $p(z | x_{obs}) \approx q_{\varphi^*}(z; x_{obs})$

Often used in a sequential manner, known as sequential NPE (SNPE), where the prior is replaced by a changing proposal distribution.



Applications of NPE+

Bayesian linear regression



Papamakarios and Murray. Fast ε-free inference of simulation models with Bayesian conditional density estimation, 2018. Wagner-Carena et al. From Images to Dark Matter: End-To-End Inference of Substructure From Hundreds of Strong Gravitational Lenses, 2023. Dingledein et al. Amortized template matching of molecular conformations from cryoelectron microscopy images using simulation-based inference, 2025.

Astrophysics



Recall the basics of VI:

- **Step 1:** Specify a *divergence* $\mathscr{D}(q; \pi)$ between variational density q(z) and target $\pi(z)$
- **Step 2:** Find the member of the family \hat{Q} that minimizes the divergence, i.e.,
 - $q^* = \arg\min \mathcal{D}(q; \pi)$ $q \in Q$
- **Step 3:** Use variational density q^* instead of target $\pi(z)$ in downstream tasks: e.g.,
 - $\mathbb{E}_{\pi(z)}[f(z)] \approx \mathbb{E}_{a^{*}(z)}[f(z)]$





Recall the basics of VI: $\mathbb{E}_{\pi(z)}[f(z)]$

- **Step 1:** Specify a *divergence* $\mathscr{D}(q; \pi)$ between variational density q(z) and target $\pi(z)$
- **Step 2:** Find the member of the family \hat{Q} that minimizes the divergence, i.e.,
 - $q^* = \arg\min \mathcal{D}(q;\pi)$ $q \in Q$
- **Step 3:** Use variational density q^* instead of target $\pi(z)$ in downstream tasks: e.g.,

$$] \approx \mathbb{E}_{q^{*}(z)}[f(z)]$$







Recall the basics of VI: $\mathbb{E}_{\pi(z)}[f(z)]$

Expressivity

- correlations
- skew and kurtosis
- multi-modality



- **Step 1:** Specify a *divergence* $\mathscr{D}(q; \pi)$ between variational density q(z) and target $\pi(z)$
- **Step 2:** Find the member of the family Q that minimizes the divergence, i.e.,
 - $q^* = \arg\min \mathcal{D}(q; \pi)$ $q \in Q$
- **Step 3:** Use variational density q^* instead of target $\pi(z)$ in downstream tasks: e.g.,

$$] \approx \mathbb{E}_{q^{*}(z)}[f(z)]$$





Recall the basics of VI: $\mathbb{E}_{\pi(z)}[f(z)]$

Expressivity

- correlations
- skew and kurtosis
- multi-modality



- **Step 1:** Specify a *divergence* $\mathscr{D}(q; \pi)$ between variational density q(z) and target $\pi(z)$
- **Step 2:** Find the member of the family Q that minimizes the divergence, i.e.,
 - $q^* = \arg\min \mathcal{D}(q; \pi)$ $q \in Q$
- **Step 3:** Use variational density q^* instead of target $\pi(z)$ in downstream tasks: e.g.,

$$] \approx \mathbb{E}_{q^{*}(z)}[f(z)]$$

- Easy and fast to evaluate q(z)
- Fast sampling from its distribution (e.g., to approximate expectations)
- Ease of optimization





Expressivity

- correlations
- skew and kurtosis
- multi-modality

- Easy and fast to evaluate q(z)
- Efficient sampling from its distribution (e.g., to approximate expectations)
- Ease of optimization



Expressivity

- correlations
- skew and kurtosis
- multi-modality

Target distribution



- Easy and fast to evaluate q(z)
- Efficient sampling from its distribution (e.g., to approximate expectations)
- Ease of optimization



Expressivity

- correlations
- skew and kurtosis
- multi-modality

Target distribution



Factorized Gaussian



- Easy and fast to evaluate q(z)
- Efficient sampling from its distribution (e.g., to approximate expectations)
- Ease of optimization





Expressivity

- correlations
- skew and kurtosis
- multi-modality

Target distribution



Factorized Gaussian



Tractability

- Easy and fast to evaluate q(z)
- Efficient sampling from its distribution (e.g., to approximate expectations)
- Ease of optimization



Energy-based model



 $q(z) = \frac{1}{Z_{\phi}} e^{-E_{\phi}(z)}$



Part III: Advances Score matching, flow matching, and diffusion



Black-box inference: no model-specific derivations, general targets, accessible to scientists

Black-box inference: no model-specific derivations, general targets, accessible to scientists

Variational inference (VI): casts the inference problem as an optimization problem: find the member of a simpler family that is closest to the target distribution

Black-box inference: no model-specific derivations, general targets, accessible to scientists

Variational inference (VI): casts the inference problem as an optimization problem: find the member of a simpler family that is closest to the target distribution



Black-box inference: no model-specific derivations, general targets, accessible to scientists

Variational inference (VI): casts the inference problem as an optimization problem: find the member of a simpler family that is closest to the target distribution

Expressivity

- correlations
- skew and kurtosis
- multi-modality



Black-box inference: no model-specific derivations, general targets, accessible to scientists

Variational inference (VI): casts the inference problem as an optimization problem: find the member of a simpler family that is closest to the target distribution

Expressivity

- correlations
- skew and kurtosis
- multi-modality



Tractability

- Easy to evaluate q(z)
- Can sample from its distribution (e.g., to approximate expectations)
- Ease of optimization

Black-box inference: no model-specific derivations, general targets, accessible to scientists

Variational inference (VI): casts the inference problem as an optimization problem: find the member of a simpler family that is closest to the target distribution

Expressivity

- correlations
- skew and kurtosis
- multi-modality



Tractability

- Easy to evaluate q(z)
- Can sample from its distribution (e.g., to approximate expectations)
- Ease of optimization

Most common in VI: factorized (mean-field) Gaussian — not so expressive but very tractable
Black-box inference, variational inference, and misspecification

Black-box inference: no model-specific derivations, general targets, accessible to scientists

Variational inference (VI): casts the inference problem as an optimization problem: find the member of a simpler family that is closest to the target distribution

Expressivity

- correlations
- skew and kurtosis
- multi-modality



Most common in VI: factorized (mean-field) Gaussian — not so expressive but very tractable

Many successes with this family for BBVI via stochastic gradient methods

Tractability

- Easy to evaluate q(z)
- Can sample from its distribution (e.g., to approximate expectations)
- Ease of optimization

57

More expressive families face challenges in optimization: can have very slow convergence due to noisy gradients, sensitivity to learning rates, batch sizes, etc.



More expressive families face challenges in optimization: can have very slow convergence due to noisy gradients, sensitivity to learning rates, batch sizes, etc.

Happens even in Gaussian (with full covariance)!



More expressive families face challenges in optimization: can have very slow convergence due to noisy gradients, sensitivity to learning rates, batch sizes, etc.

capitalize on the *structure* of the variational family? Structure = Scores

- Happens even in Gaussian (with full covariance)!
- **Question:** What if we keep the "black box" nature of the target distribution, but



More expressive families face challenges in optimization: can have very slow convergence due to noisy gradients, sensitivity to learning rates, batch sizes, etc.

capitalize on the *structure* of the variational family? Structure = Scores



- Happens even in Gaussian (with full covariance)!
- **Question:** What if we keep the "black box" nature of the target distribution, but



Why scores? "Black box": don't need $\pi(z)$, just its ("Stein") score $\nabla_z \log \pi(z)$ (Avoids normalizing constants)



Why scores? "Black box": don't need $\pi(z)$, just its ("Stein") score $\nabla_z \log \pi(z)$ (Avoids normalizing constants)

Given samples, how to make the scores close?



Why scores? **"Black box"**: don't need $\pi(z)$, just its ("Stein") score $\nabla_z \log \pi(z)$ (Avoids normalizing constants) Given samples, how to make the scores close?



Latent space



Why scores? **"Black box"**: don't need $\pi(z)$, just its ("Stein") score $\nabla_z \log \pi(z)$ (Avoids normalizing constants) Given samples, how to make the scores close?





Why scores? **"Black box"**: don't need $\pi(z)$, just its ("Stein") score $\nabla_z \log \pi(z)$ (Avoids normalizing constants) Given samples, how to make the scores close?





Why scores? "Black box": don't need $\pi(z)$, just its ("Stein") score $\nabla_{z} \log \pi(z)$ (Avoids normalizing constants) Given samples, how to make the scores close?



- $\log q(z)$ is a quadratic
- $\nabla \log q(z) = \Sigma_q^{-1}(z \mu_q)$





(Avoids normalizing constants) Given samples, how to make the scores close?



Why scores? "Black box": don't need $\pi(z)$, just its ("Stein") score $\nabla_{\tau} \log \pi(z)$

- $\log q(z)$ is a quadratic
- $\nabla \log q(z) = \Sigma_q^{-1}(z \mu_q)$

Latent space





(Avoids normalizing constants) Given samples, how to make the scores close?



Why scores? "Black box": don't need $\pi(z)$, just its ("Stein") score $\nabla_{z} \log \pi(z)$

- $\log q(z)$ is a quadratic
- $\nabla \log q(z) = \Sigma_q^{-1}(z \mu_q)$

Latent space





(Avoids normalizing constants) Given samples, how to make the scores close?



Why scores? "Black box": don't need $\pi(z)$, just its ("Stein") score $\nabla_{z} \log \pi(z)$

- $\log q(z)$ is a quadratic
- $\nabla \log q(z) = \Sigma_q^{-1}(z \mu_q)$





(Avoids normalizing constants) Given samples, how to make the scores close?



Why scores? "Black box": don't need $\pi(z)$, just its ("Stein") score $\nabla_{z} \log \pi(z)$

For Gaussian Q:

- $\log q(z)$ is a quadratic
- $\nabla \log q(z) = \Sigma_q^{-1}(z \mu_q)$

Under mild conditions, $\log p = \log q \iff \nabla \log p = \nabla \log q$





(Avoids normalizing constants) Given samples, how to make the scores close?



Why scores? "Black box": don't need $\pi(z)$, just its ("Stein") score $\nabla_{7}\log \pi(z)$

For Gaussian Q:

- $\log q(z)$ is a quadratic
- $\nabla \log q(z) = \Sigma_q^{-1}(z \mu_q)$

Under mild conditions, $\log p = \log q \iff \nabla \log p = \nabla \log q$

But: can't match scores exactly at multiple points for non-Gaussian targets





(Avoids normalizing constants) Given samples, how to make the scores close?



Modi et al.. Variational inference with Gaussian score matching. NeurIPS 2023.

Why scores? "Black box": don't need $\pi(z)$, just its ("Stein") score $\nabla_{7}\log \pi(z)$

For Gaussian Q:

- $\log q(z)$ is a quadratic
- $\nabla \log q(z) = \Sigma_q^{-1}(z \mu_q)$

Under mild conditions, $\log p = \log q \iff \nabla \log p = \nabla \log q$

But: can't match scores exactly at multiple points for non-Gaussian targets

Instead: Formulate variational inference with a score-based divergence





Fisher Divergence: measures the agreement between the scores of q and π : $\mathscr{D}(q;\pi) = \int \| \nabla_z \log q$

Hyvarinen. Estimation of Non-Normalized Statistical Models by Score Matching. JMLR, 2005.

$$q(z) - \nabla_z \log \pi(z) \parallel^2 q(z) \, dz$$



$$\mathscr{D}(q;\pi) = \int \left\| \nabla_z \log q(z) - \nabla_z \log \pi(z) \right\|^2 q(z) \, dz$$



Hyvarinen. Estimation of Non-Normalized Statistical Models by Score Matching. JMLR, 2005. Yang et al. Variational approximations using Fisher divergence. arXiv, 2019. Chen et al. Weighted Fisher divergence for high-dimensional Gaussian variational inference. arXiv, 2025



$$\mathscr{D}(q;\pi) = \int \left\| \nabla_z \log q(z) - \nabla_z \log \pi(z) \right\|^2 q(z) \, dz$$



Hyvarinen. Estimation of Non-Normalized Statistical Models by Score Matching. JMLR, 2005. Yang et al. Variational approximations using Fisher divergence. arXiv, 2019. Chen et al. Weighted Fisher divergence for high-dimensional Gaussian variational inference. arXiv, 2025



$$\mathscr{D}(q;\pi) = \int \left\| \nabla_z \log q(z) - \nabla_z \log \pi(z) \right\|^2 q(z) \, dz$$



Hyvarinen. Estimation of Non-Normalized Statistical Models by Score Matching. JMLR, 2005. Yang et al. Variational approximations using Fisher divergence. arXiv, 2019. Chen et al. Weighted Fisher divergence for high-dimensional Gaussian variational inference. arXiv, 2025



$$\mathscr{D}(q;\pi) = \int \left\| \nabla_z \log q(z) - \nabla_z \log \pi(z) \right\|^2 q(z) \, dz$$



Hyvarinen. Estimation of Non-Normalized Statistical Models by Score Matching. JMLR, 2005. Yang et al. Variational approximations using Fisher divergence. arXiv, 2019. Chen et al. Weighted Fisher divergence for high-dimensional Gaussian variational inference. arXiv, 2025



$$\mathscr{D}(q;\pi) = \int \left\| \nabla_z \log q(z) - \nabla_z \log \pi(z) \right\|^2 q(z) \, dz$$



How do we optimize this divergence? "Reparameterization trick"+SGD? Same issues as ADVI.

Hyvarinen. Estimation of Non-Normalized Statistical Models by Score Matching. JMLR, 2005. Yang et al. Variational approximations using Fisher divergence. arXiv, 2019. Chen et al. Weighted Fisher divergence for high-dimensional Gaussian variational inference. arXiv, 2025



Weighted Fisher divergence & full-covariance Gaussians

Divergence: measure the agreement between the scores of q and π :

$$\mathscr{D}(q;\pi) = \int \| \nabla_z \log q d$$

 $q(z) - \nabla_z \log \pi(z) \left\| \int_{Cov(a)}^2 q(z) dz \right\|$ $\|x\|_{\Sigma}^2 := x^{\mathsf{T}} \Sigma x$





Weighted Fisher divergence & full-covariance Gaussians

Divergence: measure the agreement between the scores of q and π :

$$\mathcal{D}(q;\pi) = \int \left\| \nabla_z \log q(z) - \nabla_z \log \pi(z) \right\|_{\operatorname{Cov}(q)}^2 q(z) \, dz$$
$$\|x\|_{\Sigma}^2 := x^{\mathsf{T}} \Sigma x$$

Properties:

- Non-negativity ($\mathscr{D} \geq 0$ with equality iff p = q)
- Invariant to affine transformations \implies leads to closed-form solutions (not true for Fisher div.)





Weighted Fisher divergence & full-covariance Gaussians

Divergence: measure the agreement between the scores of q and π :

$$\mathcal{D}(q;\pi) = \int \left\| \nabla_z \log q(z) - \nabla_z \log \pi(z) \right\|^2_{\operatorname{Cov}(q)} q(z) \, dz$$
$$\|x\|_{\Sigma}^2 := x^{\mathsf{T}} \Sigma x$$

Properties:

- Non-negativity ($\mathscr{D} \geq 0$ with equality iff p = q)
- Invariant to affine transformations \implies leads to closed-form solutions (not true for Fisher div.)

How do we optimize this divergence? Again could do "reparameterization trick" +SGD, but same issues as ADVI.





Observation for Gaussians: scores contain a lot of structure

Covariance-weighted Fisher divergence:

$$\mathscr{D}(q;\pi) = \int \left\| \nabla_z \right\|_{z}$$

$\sum_{z} \log q(z) - \nabla_{z} \log \pi(z) \| \sum_{\operatorname{Cov}(q)}^{2} q(z) \, dz$





$$\mathscr{D}(q;\pi) = \int \| \nabla_z dz$$

Form an empirical estimate:

Sample z^1, \ldots, z^B from a distribution ν (that is not q)

- $\sum_{z} \log q(z) \nabla_{z} \log \pi(z) \left\| \sum_{\operatorname{Cov}(q)}^{2} q(z) dz \right\|$





$$\mathscr{D}(q;\pi) = \int \left\| \nabla_z \log q(z) - \nabla_z \log \pi(z) \right\|_{\operatorname{Cov}(q)}^2 q(z) \, dz$$

Form an empirical estimate:

Sample z^1, \ldots, z^B from a distribution ν (that is not q)

$$\widehat{\mathcal{D}}(q;\pi) = \frac{1}{B} \sum_{b=1}^{B}$$

$$\left\| \nabla_z \log q(z^b) - \nabla_z \log \pi(z^b) \right\|_{Cov(q)}^2$$





$$\mathscr{D}(q;\pi) = \int \left\| \nabla_z \log q(z) - \nabla_z \log \pi(z) \right\|_{\operatorname{Cov}(q)}^2 q(z) \, dz$$

Form an empirical estimate: Sample z^1, \ldots, z^B from a distribution ν (that is not q)

$$\widehat{\mathcal{D}}(q;\pi) = \frac{1}{B} \sum_{b=1}^{B}$$

If π is Gaussian and B > D,

$$\left\| \nabla_z \log q(z^b) - \nabla_z \log \pi(z^b) \right\|_{Cov(q)}^2$$





$$\mathscr{D}(q;\pi) = \int \left\| \nabla_z \log q(z) - \nabla_z \log \pi(z) \right\|_{\operatorname{Cov}(q)}^2 q(z) \, dz$$

Form an empirical estimate: Sample z^1, \ldots, z^B from a distribution ν (that is not q)

$$\widehat{\mathcal{D}}(q;\pi) = \frac{1}{B} \sum_{b=1}^{B}$$

If π is Gaussian and B > D, $\hat{\mu}, \hat{\Sigma} = \arg\min\widehat{\mathscr{D}}(q_{\mu,\Sigma}; \pi)$ recovers the exact mean and covariance!

$$\left\| \nabla_z \log q(z^b) - \nabla_z \log \pi(z^b) \right\|_{Cov(q)}^2$$





$$\mathscr{D}(q;\pi) = \int \left\| \nabla_z \log q(z) - \nabla_z \log \pi(z) \right\|_{\operatorname{Cov}(q)}^2 q(z) \, dz$$

Form an empirical estimate: Sample z^1, \ldots, z^B from a distribution ν (that is not q)

$$\widehat{\mathcal{D}}(q;\pi) = \frac{1}{B} \sum_{b=1}^{B}$$

If π is Gaussian and B > D, $\hat{\mu}, \hat{\Sigma} = \arg\min \widehat{\mathscr{D}}(q_{\mu,\Sigma}; \pi)$ recovers the exact mean and covariance! μ, Σ

$$\left\| \nabla_z \log q(z^b) - \nabla_z \log \pi(z^b) \right\|_{\operatorname{Cov}(q)}^2$$

- E.g., for D = 1, only need 2 points! Compare to: 2 samples has error $\mathcal{O}(1/\sqrt{2})$





$$\mathscr{D}(q;\pi) = \int \| \nabla_z$$

Form an empirical estimate: Sample z^1, \ldots, z^B from a distribution ν (that is not q)

$$\widehat{\mathcal{D}}(q;\pi) = \frac{1}{B} \sum_{b=1}^{B}$$

If π is Gaussian and B > D, $\hat{\mu}, \hat{\Sigma} = \arg\min \widehat{\mathscr{D}}(q_{\mu,\Sigma}; \pi)$ recovers the exact mean and covariance! μ, Σ

How to turn this into an algorithm? Need a sampling distribution.

Observation for Gaussians: scores contain a lot of structure

- $\sum_{z} \log q(z) \nabla_{z} \log \pi(z) \left\| \sum_{\operatorname{Cov}(q)}^{2} q(z) \, dz \right\|_{\operatorname{Cov}(q)}$

$$\left\| \nabla_z \log q(z^b) - \nabla_z \log \pi(z^b) \right\|_{Cov(q)}^2$$

E.g., for D = 1, only need 2 points! Compare to: 2 samples has error $\mathcal{O}(1/\sqrt{2})$





Score-based variational inference

Monte Carlo estimate (unbiased): **Score-based VI:** Want $q^* = \arg \min \mathcal{D}(q; \pi)$ Sample z^1, \ldots, z^B from q $q \in Q$ $\nabla_z \log q(z^b) - \nabla_z \log \pi(z^b) \|_{Cov(q)}^2$

$$\widehat{\mathscr{D}}(q;\pi) = \frac{1}{B} \sum_{b=1}^{B} \| \nabla_{a}$$



Score-based variational inference

Score-based VI: Want $q^* = \arg \min \mathcal{D}(q; \pi)$ $q \in Q$

$$\widehat{\mathcal{D}}(q;\pi) = \frac{1}{B} \sum_{b=1}^{B} \| \nabla_{z}$$

But: cannot simultaneously sample from and optimize over q!

Monte Carlo estimate (unbiased): Sample z^1, \ldots, z^B from q $v_z \log q(z^b) - \nabla_z \log \pi(z^b) \|_{Cov(q)}^2$



Score-based variational inference

Score-based VI: Want $q^* = \arg \min \mathcal{D}(q; \pi)$ $q \in Q$

$$\widehat{\mathcal{D}}(q;\pi) = \frac{1}{B} \sum_{b=1}^{B} \| \nabla_{z}$$

But: cannot simultaneously sample from and optimize over q!

Iterative procedure:



Monte Carlo estimate (unbiased): Sample z^1, \ldots, z^B from q $\sqrt{1}_{z}\log q(z^{b}) - \nabla_{z}\log \pi(z^{b}) \|_{Cov(q)}^{2}$




Score-based VI: Want $q^* = \arg \min \mathcal{D}(q; \pi)$ $q \in Q$

$$\widehat{\mathcal{D}}(q;\pi) = \frac{1}{B} \sum_{b=1}^{B} \| \nabla_{z}$$

But: cannot simultaneously sample from and optimize over q!

Iterative procedure:







Score-based VI: Want $q^* = \arg \min \mathcal{D}(q; \pi)$ $q \in Q$

$$\widehat{\mathcal{D}}(q;\pi) = \frac{1}{B} \sum_{b=1}^{B} \| \nabla_{z}$$

But: cannot simultaneously sample from and optimize over q!

Iterative procedure:







Score-based VI: Want $q^* = \arg \min \mathcal{D}(q; \pi)$ $q \in Q$

$$\widehat{\mathcal{D}}(q;\pi) = \frac{1}{B} \sum_{b=1}^{B} \| \nabla_{z}$$

But: cannot simultaneously sample from and optimize over q!

Iterative procedure:





Score-based VI: Want $q^* = \arg \min \mathcal{D}(q; \pi)$ $q \in Q$

$$\widehat{\mathcal{D}}(q;\pi) = \frac{1}{B} \sum_{b=1}^{B} \| \nabla_{z}$$

But: cannot simultaneously sample from and optimize over q!

Iterative procedure:





Score-based VI: Want $q^* = \arg \min \mathcal{D}(q; \pi)$ $q \in Q$

$$\widehat{\mathcal{D}}(q;\pi) = \frac{1}{B} \sum_{b=1}^{B} \| \nabla_{z}$$

But: cannot simultaneously sample from and optimize over q!

Iterative procedure:





Score-based VI: Want $q^* = \arg \min \mathcal{D}(q; \pi)$ $q \in Q$

$$\widehat{\mathcal{D}}(q;\pi) = \frac{1}{B} \sum_{b=1}^{B} \| \nabla_{z}$$

But: cannot simultaneously sample from and optimize over q!

Iterative procedure:





Score-based VI: Want $q^* = \arg \min \mathcal{D}(q; \pi)$ $q \in Q$

$$\widehat{\mathcal{D}}(q;\pi) = \frac{1}{B} \sum_{b=1}^{B} \| \nabla_{z}$$

But: cannot simultaneously sample from and optimize over q!

Iterative procedu

 π

re:
$$q^*$$
 q_{t+1}
small

Introduces bias: regularize the score divergence so we don't move too far from q_t





Score-based VI: Want $q^* = \arg \min \mathcal{D}(q; \pi)$ $q \in Q$

$$\widehat{\mathcal{D}}(q;\pi) = \frac{1}{B} \sum_{b=1}^{B} \| \nabla_{z}$$

But: cannot simultaneously sample from and optimize over q!

 π

Iterative procedu

re:
$$q^*$$
 q_{t+1}
small

Introduces bias: regularize the score divergence so we don't move too far from q_{t}





Score-based VI: Want $q^* = \arg \min \mathcal{D}(q; \pi)$ $q \in Q$

$$\widehat{\mathcal{D}}(q;\pi) = \frac{1}{B} \sum_{b=1}^{B} \| \nabla_{z}$$

But: cannot simultaneously sample from and optimize over q!

Iterative procedure:



Introduces bias: regularize the score divergence so we don't move too far from q_{t}

Proximal point algorithm:

Here, the *global minimum* of the subproblem can be computed analytically!

Monte Carlo estimate (unbiased): Sample z^1, \ldots, z^B from q $\sum_{z} \log q(z^b) - \nabla_z \log \pi(z^b) \|_{Cov(q)}^2$

 $q^* \qquad q_{t+1} = \arg\min_{q \in \mathcal{Q}} \widehat{\mathcal{Q}}_{q_t}(q; \pi) + \operatorname{regularization}(q_t; q)$ $\operatorname{small step} q_t \qquad q_0 \quad Q$





Iterate:

Batch step: given current variational estimate q_t Sample a batch of points z^1, \ldots, z^B (call B the batch size) Match step:

- Compute *empirical* score-based divergence $\widehat{\mathcal{D}}_{q_t}(q; \pi)$ on batch of points
- Fit new variational distribution $q_{t+1} = \arg \min_{q \in Q} \widehat{\mathcal{D}}_{q_t}(q; \pi) + \frac{2}{\lambda_t} \operatorname{KL}(q_t; q)$





Iterate:

Batch step: given current variational estimate q_t Sample a batch of points z^1, \ldots, z^B (call B the batch size) Match step:

Compute *empirical* score-based divergence $\widehat{\mathcal{D}}_{q_t}(q; \pi)$ on batch of points Fit new variational distribution $q_{t+1} = \arg \min_{q \in Q} \widehat{\mathcal{D}}_{q_t}(q; \pi) + \frac{2}{\lambda_t} \operatorname{KL}(q_t; q)$ empirical divergence





Iterate:

Batch step: given current variation Sample a *batch* of points z^1 , . Compute *empirical* score-base Match step:

Fit new variational distribution

hal estimate
$$q_t$$

..., z^B (call B the batch size)
ed divergence $\widehat{\mathscr{D}}_{q_t}(q; \pi)$ on batch of points
n $q_{t+1} = \arg\min_{q \in \mathscr{Q}} \widehat{\mathscr{D}}_{q_t}(q; \pi) + \frac{2}{\lambda_t} \operatorname{KL}(q_t; q)$
empirical "learning
divergence rate"





Iterate:

Batch step: given current variation Sample a *batch* of points z^1 , . Compute *empirical* score-base Match step:

Fit new variational distribution

hal estimate
$$q_t$$

..., z^B (call B the batch size)
ed divergence $\widehat{\mathscr{D}}_{q_t}(q; \pi)$ on batch of points
n $q_{t+1} = \arg\min_{q \in \mathscr{Q}} \widehat{\mathscr{D}}_{q_t}(q; \pi) + \frac{2}{\lambda_t} \operatorname{KL}(q_t; q)$
empirical "learning regularizer
divergence rate"





Iterate:

Batch step: given current variation Sample a *batch* of points z^1 , . Compute *empirical* score-base Match step: Fit new variational distribution

In this problem, the global minimum of the subproblem can be computed analytically!

hal estimate
$$q_t$$

..., z^B (call B the batch size)
ed divergence $\widehat{\mathscr{D}}_{q_t}(q; \pi)$ on batch of points
in $q_{t+1} = \arg\min_{q \in \mathscr{Q}} \widehat{\mathscr{D}}_{q_t}(q; \pi) + \frac{2}{\lambda_t} \operatorname{KL}(q_t; q)$
empirical "learning regularizer
divergence rate"





Iterate:

Batch step: given current variation Sample a *batch* of points z^1 , . Compute *empirical* score-base Match step:

Fit new variational distribution

In this problem, the *global minimum* of the subproblem can be computed analytically! Σ_{t+1} update: solve a quadratic matrix equation

$$\Sigma_{t+1} U_t \Sigma_{t+1} + \Sigma_{t+1} = V_t$$

hal estimate
$$q_t$$

..., z^B (call B the batch size)
ed divergence $\widehat{\mathscr{D}}_{q_t}(q; \pi)$ on batch of points
in $q_{t+1} = \arg\min_{q \in \mathscr{Q}} \widehat{\mathscr{D}}_{q_t}(q; \pi) + \frac{2}{\lambda_t} \operatorname{KL}(q_t; q)$
empirical "learning regularizer
divergence rate"





Iterate:

Batch step: given current variation Sample a *batch* of points z^1 , . Compute *empirical* score-base Match step:

Fit new variational distribution

In this problem, the *global minimum* of the subproblem can be computed analytically! Σ_{t+1} update: solve a quadratic matrix equation

$$\Sigma_{t+1} U_t \Sigma_{t+1} + \Sigma_{t+1} = V_t$$

Matrices that depend on current covariance, statistics of the batch and scores

hal estimate
$$q_t$$

..., z^B (call B the batch size)
ed divergence $\widehat{\mathscr{D}}_{q_t}(q; \pi)$ on batch of points
in $q_{t+1} = \arg\min_{q \in \mathscr{Q}} \widehat{\mathscr{D}}_{q_t}(q; \pi) + \frac{2}{\lambda_t} \operatorname{KL}(q_t; q)$
empirical "learning regularizer
divergence rate"





Iterate:

Batch step: given current variation Sample a *batch* of points z^1 , . Compute *empirical* score-base Match step:

Fit new variational distribution

In this problem, the *global minimum* of the subproblem can be computed analytically!

 Σ_{t+1} update: solve a *quadratic matrix equation* mean update:

$$\Sigma_{t+1} U_t \Sigma_{t+1} + \Sigma_{t+1} = V_t$$

Matrices that depend on current covariance, statistics of the batch and scores

hal estimate
$$q_t$$

..., z^B (call B the batch size)
ed divergence $\widehat{\mathscr{D}}_{q_t}(q; \pi)$ on batch of points
in $q_{t+1} = \arg\min_{q \in \mathscr{Q}} \widehat{\mathscr{D}}_{q_t}(q; \pi) + \frac{2}{\lambda_t} \operatorname{KL}(q_t; q)$
empirical "learning regularizer
divergence rate"

$$\mu_{t+1} = \mu_t + \frac{\lambda_t}{1+\lambda_t} \Sigma_{t+1} \bar{g}$$





Iterate:

Batch step: given current variation Sample a *batch* of points z^1 , . Compute *empirical* score-base Match step:

Fit new variational distribution

In this problem, the *global minimum* of the subproblem can be computed analytically!

 Σ_{t+1} update: solve a quadratic matrix equation

$$\Sigma_{t+1} U_t \Sigma_{t+1} + \Sigma_{t+1} = V_t$$

Matrices that depend on current covariance, statistics of the batch and scores

hal estimate
$$q_t$$

..., z^B (call B the batch size)
ed divergence $\widehat{\mathscr{D}}_{q_t}(q; \pi)$ on batch of points
in $q_{t+1} = \arg\min_{q \in \mathscr{Q}} \widehat{\mathscr{D}}_{q_t}(q; \pi) + \frac{2}{\lambda_t} \operatorname{KL}(q_t; q)$
empirical "learning regularizer
divergence rate"

mean update:

$$\mu_{t+1} = \mu_t + \frac{\lambda_t}{1+\lambda_t} \Sigma_{t+1} \bar{g}$$

mean of scores computed on the batch





Application: deep generative modeling

Data: high-dimensional object (image) $x_n \in \mathbb{R}^{3072}$

Lower-dimensional latent representation $z_n \in \mathbb{R}^{256}$

Cai et al. Batch and match: black-box variational inference with a score-based divergence. ICML 2024.





Application: deep generative modeling Deep generative model: Deep generative model:

Data: high-dimensional object (image) $x_n \in \mathbb{R}^{3072}$

Lower-dimensional latent representation $z_n \in \mathbb{R}^{256}$

Cai et al. Batch and match: black-box variational inference with a score-based divergence. ICML 2024.



Application: deep generative modeling **Deep generative model:** Data: high-dimensional object (image) $x_n \in \mathbb{R}^{3072}$ $z_n \sim \mathcal{N}(0,I)$ $x_n | z_n \sim \mathcal{N}(\Omega(z_n, \hat{\theta}), \sigma^2 I)$

Lower-dimensional latent representation $z_n \in \mathbb{R}^{256}$

Cai et al. Batch and match: black-box variational inference with a score-based divergence. ICML 2024.





Application: deep generative modeling

Data: high-dimensional object (image) $x_n \in \mathbb{R}^{3072}$

Lower-dimensional latent representation $z_n \in \mathbb{R}^{256}$

Cai et al. Batch and match: black-box variational inference with a score-based divergence. ICML 2024.



Deep generative model:

 $z_n \sim \mathcal{N}(0,I)$

 $x_n | z_n \sim \mathcal{N}(\Omega(z_n, \hat{\theta}), \sigma^2 I)$

Parameters of the neural network (pre-trained to maximize likelihood)





Problem: given a test image x', approximate p(z' | x') using VI

Cai et al. Batch and match: black-box variational inference with a score-based divergence. ICML 2024.

Parameters of the neural network (pre-trained to maximize likelihood)





Application: deep generative modeling **Deep generative model:** Data: high-dimensional object (image) $x_n \in \mathbb{R}^{3072}$ $z_n \sim \mathcal{N}(0,I)$ $x_n | z_n \sim \mathcal{N}(\Omega(z_n, \hat{\theta}), \sigma^2 I)$

Lower-dimensional latent representation $z_n \in \mathbb{R}^{256}$

Problem: given a test image x', approximate p(z' | x') using VI

Reconstruct x' by feeding the posterior mean estimate into the neural network.

Cai et al. Batch and match: black-box variational inference with a score-based divergence. ICML 2024.

Parameters of the neural network (pre-trained to maximize likelihood)





Problem: given a test image x', approximate p(z' | x') using VI

Reconstruct x' by feeding the posterior mean estimate into the neural network.



Cai et al. Batch and match: black-box variational inference with a score-based divergence. ICML 2024.

Parameters of the neural network (pre-trained to maximize likelihood)





1000





Problem: given a test image x', approximate p(z' | x') using VI

Reconstruct x' by feeding the posterior mean estimate into the neural network.



Cai et al. Batch and match: black-box variational inference with a score-based divergence. ICML 2024.

Parameters of the neural network (pre-trained to maximize likelihood)

B=300 Vean Squared Error 0 1:0 0 5:0 Amortized 1000 VI 200 300 100 0 Wall time (s)

Image reconstruction error





Problem: given a test image x', approximate p(z' | x') using VI

Reconstruct x' by feeding the posterior mean estimate into the neural network.



Cai et al. Batch and match: black-box variational inference with a score-based divergence. ICML 2024.

Parameters of the neural network (pre-trained to maximize likelihood)

B=300 Vean Squared Error 0 1:0 0 5:0 Amortized 1000 VI 200 300 100 Wall time (s) BaM with B = 300 converges faster than ADVI with any batch size.

Image reconstruction error





Cai, Modi, Margossian, Gower, Blei, Saul. EigenVI: score-based variational inference with orthogonal function expansions. *NeurIPS* 2024. Modi & Cai, Saul. Batch, match, and patch: low-rank approximations for score-based variational inference. *AISTATS* 2025.



- A "patch" step: we project the covariance update to low-rank + diagonal
- The resulting covariance update has $\mathcal{O}(D)$ cost

Cai, Modi, Margossian, Gower, Blei, Saul. EigenVI: score-based variational inference with orthogonal function expansions. *NeurIPS* 2024. Modi & Cai, Saul. Batch, match, and patch: low-rank approximations for score-based variational inference. AISTATS 2025.



- A "patch" step: we project the covariance update to low-rank + diagonal
- The resulting covariance update has $\mathcal{O}(D)$ cost

How can we use score-based VI for Bayesian inference with massive data sets?

• Can use BaM with a noisy score $\hat{\nabla}\log p(z)$, computed from a subsampled data set

Cai, Modi, Margossian, Gower, Blei, Saul. EigenVI: score-based variational inference with orthogonal function expansions. *NeurIPS* 2024. Modi & Cai, Saul. Batch, match, and patch: low-rank approximations for score-based variational inference. AISTATS 2025.



- A "patch" step: we project the covariance update to low-rank + diagonal
- The resulting covariance update has $\mathcal{O}(D)$ cost

How can we use score-based VI for Bayesian inference with massive data sets?

• Can use BaM with a noisy score $\hat{\nabla}\log p(z)$, computed from a subsampled data set

What about richer variational families?

Cai, Modi, Margossian, Gower, Blei, Saul. EigenVI: score-based variational inference with orthogonal function expansions. *NeurIPS* 2024. Modi & Cai, Saul. Batch, match, and patch: low-rank approximations for score-based variational inference. AISTATS 2025.



- A "patch" step: we project the covariance update to low-rank + diagonal
- The resulting covariance update has $\mathcal{O}(D)$ cost

How can we use score-based VI for Bayesian inference with massive data sets?

• Can use BaM with a noisy score $\hat{\nabla}\log p(z)$, computed from a subsampled data set

What about richer variational families?

• EigenVI: orthogonal function expansions

Cai, Modi, Margossian, Gower, Blei, Saul. EigenVI: score-based variational inference with orthogonal function expansions. *NeurIPS* 2024. Modi & Cai, Saul. Batch, match, and patch: low-rank approximations for score-based variational inference. AISTATS 2025.



- A "patch" step: we project the covariance update to low-rank + diagonal
- The resulting covariance update has $\mathcal{O}(D)$ cost

How can we use score-based VI for Bayesian inference with massive data sets?

• Can use BaM with a noisy score $\hat{\nabla}\log p(z)$, computed from a subsampled data set

What about richer variational families?

EigenVI: orthogonal function expansions

Cai, Modi, Margossian, Gower, Blei, Saul. EigenVI: score-based variational inference with orthogonal function expansions. NeurIPS 2024. Modi & Cai, Saul. Batch, match, and patch: low-rank approximations for score-based variational inference. AISTATS 2025.





- A "patch" step: we project the covariance update to low-rank + diagonal
- The resulting covariance update has $\mathcal{O}(D)$ cost

How can we use score-based VI for Bayesian inference with massive data sets?

• Can use BaM with a noisy score $\hat{\nabla}\log p(z)$, computed from a subsampled data set

What about richer variational families?

EigenVI: orthogonal function expansions

Cai, Modi, Margossian, Gower, Blei, Saul. EigenVI: score-based variational inference with orthogonal function expansions. NeurIPS 2024. Modi & Cai, Saul. Batch, match, and patch: low-rank approximations for score-based variational inference. AISTATS 2025.





- A "patch" step: we project the covariance update to low-rank + diagonal
- The resulting covariance update has $\mathcal{O}(D)$ cost

How can we use score-based VI for Bayesian inference with massive data sets?

• Can use BaM with a noisy score $\hat{\nabla}\log p(z)$, computed from a subsampled data set

What about richer variational families?

EigenVI: orthogonal function expansions

Cai, Modi, Margossian, Gower, Blei, Saul. EigenVI: score-based variational inference with orthogonal function expansions. NeurIPS 2024. Modi & Cai, Saul. Batch, match, and patch: low-rank approximations for score-based variational inference. AISTATS 2025.





- A "patch" step: we project the covariance update to low-rank + diagonal
- The resulting covariance update has $\mathcal{O}(D)$ cost

How can we use score-based VI for Bayesian inference with massive data sets?

• Can use BaM with a noisy score $\hat{\nabla}\log p(z)$, computed from a subsampled data set

What about richer variational families?

EigenVI: orthogonal function expansions

Cai, Modi, Margossian, Gower, Blei, Saul. EigenVI: score-based variational inference with orthogonal function expansions. NeurIPS 2024. Modi & Cai, Saul. Batch, match, and patch: low-rank approximations for score-based variational inference. AISTATS 2025.




Extensions: high dimensional, large data sets

- A "patch" step: we project the covariance update to low-rank + diagonal
- The resulting covariance update has $\mathcal{O}(D)$ cost

How can we use score-based VI for Bayesian inference with massive data sets?

• Can use BaM with a noisy score $\hat{\nabla}\log p(z)$, computed from a subsampled data set

What about richer variational families?

EigenVI: orthogonal function expansions

Cai, Modi, Margossian, Gower, Blei, Saul. EigenVI: score-based variational inference with orthogonal function expansions. NeurIPS 2024. Modi & Cai, Saul. Batch, match, and patch: low-rank approximations for score-based variational inference. AISTATS 2025.

Assumes we can fit a covariance with $\mathcal{O}(D^2)$ parameters. What about in high dimensions?





Extensions: high dimensional, large data sets

- A "patch" step: we project the covariance update to low-rank + diagonal
- The resulting covariance update has $\mathcal{O}(D)$ cost

How can we use score-based VI for Bayesian inference with massive data sets?

• Can use BaM with a noisy score $\hat{\nabla}\log p(z)$, computed from a subsampled data set



Modi & Cai, Saul. Batch, match, and patch: low-rank approximations for score-based variational inference. AISTATS 2025.

Assumes we can fit a covariance with $\mathcal{O}(D^2)$ parameters. What about in high dimensions?



Extensions: high dimensional, large data sets

- A "patch" step: we project the covariance update to low-rank + diagonal
- The resulting covariance update has $\mathcal{O}(D)$ cost

How can we use score-based VI for Bayesian inference with massive data sets?

• Can use BaM with a noisy score $\hat{\nabla}\log p(z)$, computed from a subsampled data set



Modi & Cai, Saul. Batch, match, and patch: low-rank approximations for score-based variational inference. AISTATS 2025.

Assumes we can fit a covariance with $\mathcal{O}(D^2)$ parameters. What about in high dimensions?



Summary of score-based VI

We discussed score-based variational inference

- black-box and only require computing the score of the target $V_7 \log \pi(z)$
- VI with score-based divergences, e.g., Fisher divergence
 - For Gaussian Q, can apply reparameterization trick to optimize
- VI with a covariance-weighted Fisher divergence
 - Has some nice properties, including computational benefits
 - For Gaussian Q, can also use reparameterization trick to optimize
 - Batch & match: Gaussian full covariance, closed-form proximal updates
- Fast convergence for near-Gaussian targets
- Extensions of score-based VI to high dimensions, large data sets, and non-Gaussian families

67

Designing q Distributions





Normalizing flows

Auxiliary variables & mixture distributions

Use many layers (hierarchical) + continuous-time limit

Diffusion SDE & Continuous Normalizing Flow posteriors for approximate inference

• Construct q(z|x) as a (hierarchical) mixture distribution

 $q(z|x) = \int q(z | a, x) q(a|x) da$

• *a* is the auxiliary variable used to enrich the approximate posterior

• Example: Mixture of Gaussians

 $a \sim q(a|x) = Categorical(\pi_1, ..., \pi_K)$ $z \sim q(z | a, x) = N(\theta; m_a, \Sigma_a)$

Can be very flexible with many components!



• Construct q(z|x) as a (hierarchical) mixture distribution

$$q(z|x) = \int q(z | a, x) q(a|x) da$$

- *a* is the auxiliary variable used to enrich the approximate posterior
- Now the variational lower-bound becomes intractable:

$$L(\phi) = E_{q(z|x)}[\log p(x,z)] - E_{q(z|x)}[\log q(z|x)]$$

Estimated by Monte Carlo: $a^k \sim q(a|x), z^k \sim q(z \mid a^k, x)$ Intractable density $q(z|x) = \int q(z|a, x)q(a|x) da$

• Solution: introducing an auxiliary variational lower-bound $L(\phi, r)$ with an auxiliary distribution r(a|z, x):



Agakov and Barber. An Auxiliary Variational Method. ICONIP 2004 Salimans et al. Markov Chain Monte Carlo and Variational Inference: Bridging the Gap. ICML 2015 Ranganath et al. Hierarchical Variational Models. ICML 2016

• Solution: introducing an auxiliary variational lower-bound $L(\phi, r)$ with an auxiliary distribution r(a|z, x):

Salimans et al. Markov Chain Monte Carlo and Variational Inference: Bridging the Gap. ICML 2015 Ranganath et al. Hierarchical Variational Models. ICML 2016

• Solution: introducing an auxiliary variational lower-bound $L(\phi, r)$ with an auxiliary distribution r(a|z, x):



Agakov and Barber. An Auxiliary Variational Method. ICONIP 2004 Salimans et al. Markov Chain Monte Carlo and Variational Inference: Bridging the Gap. ICML 2015 Ranganath et al. Hierarchical Variational Models. ICML 2016

• Solution: introducing an auxiliary variational lower-bound $L(\phi, r)$ with an auxiliary distribution r(a|z, x):

$$log p(x)$$

$$L(\phi) = E_{q(z|x)}[log p(x|z)] - KL[q(z|x)||p(z)]$$

$$E_{q(z|x)}[KL[q(a|z,x)||r(a|z,x)]]$$

$$L(\phi,r) = E_{q(z,a|x)}[log p(x|z)] - KL[q(z,a|x)||p(z)r(a|z,x)]$$

- Optimize r(a|z, x) to close the gap!
- $L(\phi, r)$ estimated by Monte Carlo: $a^k \sim q(a|x), z^k \sim q(z | a^k, x)$

Agakov and Barber. An Auxiliary Variational Method. ICONIP 2004 Salimans et al. Markov Chain Monte Carlo and Variational Inference: Bridging the Gap. ICML 2015 Ranganath et al. Hierarchical Variational Models. ICML 2016

• Solution: introducing an auxiliary variational lower-bound $L(\phi, r)$ with an auxiliary distribution r(a|z, x):

$$\log p(x)$$

$$L(\phi) = E_{q(z|x)}[\log p(x|z)] - KL[q(z|x)||p(z)]$$

$$L(\phi, r) = E_{q(z,a|x)}[\log p(x|z)] - KL[q(z,a|x)||p(z)r(a|z,x)]$$

$$L(\phi, r) = E_{q(z,a|x)}[\log p(x|z)] - KL[q(z,a|x)||p(z)r(a|z,x)]$$
Augmented generative model view:
$$p(x, z, a) \coloneqq p(z)p(x|z)r(a|z, x) \implies p(z, a|x) = p(z|x)r(a|z, x)$$
Agakov and Barber. An Auxiliary Variational Method. ICONIP 2004
Salimans et al. Markov Chain Monte Carlo and Variational Inference: Bridging the Gap. ICML 2015

p(

Ranganath et al. Hierarchical Variational Models. ICML 2016

- Hierarchical mixture distributions for q(z, a|x)
 - VI-MCMC hybrid: build q(z) with a Markov Chain:



Learn the transition kernel with VI:

$$z\coloneqq z_T, a=\{z_{1:T-1}\}$$

$$q(z_{1:T}|x) = \prod_{t=1}^{T} K_{\phi}(z_t|z_{t-1}), \ z_0 \coloneqq x$$

$$r(z_{1:T-1}|z_T) = \prod_{t=2}^T r(z_{t-1}|z_t)$$



Salimans et al. Markov Chain Monte Carlo and Variational Inference: Bridging the Gap. ICML 2015 Huang et al. Improving Explorability in Variational Inference with Annealed Variational Objectives. NeurIPS 2018

- Hierarchical mixture distributions for q(z, a|x)
 - VI-MCMC hybrid: build q(z) with a Markov Chain:



Learn the transition kernel with VI:

$$L(\phi, r) = E_{q(z,a|x)}[\log p(x|z)] - KL[q(z,a|x)||p(z)r(a|z,x)]$$

Equivalent to minimizing KL[q(z, a|x)||p(z, a|x)], $p(z, a|x) \coloneqq p(z|x)r(a|z, x)$:

 $\log p(x) - L(\phi, r) = KL[q(z_{1:T}|z_0) || r(z_{1:T-1}|z_T)p(z_T|x)]$



• Consider the general case $(z_T \coloneqq z)$

 $\log p(x) - L(\phi, r) = KL[q(z_{1:T}|x)||r(z_{1:T-1}|z_T)p(z_T|x)]$

```
Original model: p(x, z_T) = p(x|z_T)p(z_T)
```

```
posterior: p(z_T|x)
(the target)
```



• Consider the general case $(z_T \coloneqq z)$

 $\log p(x) - L(\phi, r) = KL[q(z_{1:T}|x) || r(z_{1:T-1}|z_T)p(z_T|x)]$

Original model: $p(x, z_T) = p(x|z_T)p(z_T)$

posterior: $p(z_T|x)$ (the target) <u>Augmented</u> model: $p(x, z_T)r(z_{1:T-1}|z_T)$ $r(z_{1:T-1}|z_T) = \prod_{t=2}^{T} r(z_{t-1}|z_t)$ approx. posterior: $q(z_{1:T}|x) = q(z_1|x) \prod_{t=2}^{T} q(z_t|z_{t-1},x)$



• Consider the general case $(z_T \coloneqq z)$

 $\log p(x) - L(\phi, r) = KL[q(z_{1:T}|x) || r(z_{1:T-1}|z_T)p(z_T|x)]$

Original model: $p(x, z_T) = p(x|z_T)p(z_T)$

posterior: $p(z_T|x)$ (the target) $\frac{\text{Gaussian (fixed)}}{\text{Augmented model:}} \quad p(x, z_T) \prod_{t=2}^T r(z_{t-1} | z_t)$ $\frac{P(x, z_T) \prod_{t=2}^T r(z_{t-1} | z_t)}{r(z_{t-1} | z_t)} = N(z^{t-1}; \sqrt{1 - \beta_t} z_t, \beta_t I)$ $p(x, z_T) \prod_{t=2}^T r(z_{t-1} | z_t)$ $p(x, z_T) \prod_{t=2}^T r(z_t, \beta_t I)$



• Consider the general case $(z_T \coloneqq z)$

 $\log p(x) - L(\phi, r) = KL[q(z_{1:T}|x) || r(z_{1:T-1}|z_T)p(z_T|x)]$



• Consider posterior approximation case:





• Re-define the time index: $t \to t/T$, $z_t \to z_{t/T}$, and take limit $T \to \infty$



He et al. No Trick, No Treat: Pursuits and Challenges Towards Simulation-free Training of Neural Samplers. arXiv:2502.06685 Vargas et al. Denoising Diffusion Samplers. ICLR 2023

• Re-define the time index: $t \to t/T$, $z_t \to z_{t/T}$, and take limit $T \to \infty$



He et al. No Trick, No Treat: Pursuits and Challenges Towards Simulation-free Training of Neural Samplers. arXiv:2502.06685 Vargas et al. Denoising Diffusion Samplers. ICLR 2023

• Re-define the time index: $t \to t/T$, $z_t \to z_{t/T}$, and take limit $T \to \infty$

 $\begin{array}{ll} \underline{\mathsf{SDE}} & r(z_{[0,1]}) \text{ satisfying } r(z_1) = p(z_1|x), \\ \underline{\mathsf{Augmented}} \ \mathsf{model} & dz_t = f_t(z_t)dt + \gamma_t dW_t, t \colon 1 \to 0 \\ & (z_1 \coloneqq z) \end{array}$

approx. posterior: $q(z_{[0,1]}|x): q(z_0|x) = N(m, \Sigma)$,

Objective:

• VI-style: $KL[q(z_{[0,1]}|x)||r(z_{[0,1]})]$

• Conditions for a well-defined KL divergence between SDEs:

 $dz_t = \mu_t(z_t, x)dt + \sigma_t d\overline{W}_t, t: 0 \to 1$

- Going in the same direction for analytic form of the KL: reverse either the r SDE or the q SDE

He et al. No Trick, No Treat: Pursuits and Challenges Towards Simulation-free Training of Neural Samplers. arXiv:2502.06685 Vargas et al. Denoising Diffusion Samplers. ICLR 2023 Li et al. Scalable Gradients for Stochastic Differential Equations. AISTATS 2020 Song et al. Score-Based Generative Modeling through Stochastic Differential Equations. ICLR 2021

Continuity Equation

• Forward SDE evolution of z_t , $t: 0 \rightarrow 1$

$$dz_t = \underline{f_t(z_t)}dt + \sigma_t(z_t)dW_t, \quad z_0 \sim p_0(z_0)$$

("forward drift")

• Fokker-Planck Equation: marginal density evolution of $p_t(z), t: 0 \to 1$, for any z (not necessarily for $z = z_t$) $\partial_t p_t(z) = -\nabla_z \cdot \left(f_t(z) p_t(z) \right) + \frac{1}{2} \nabla_z^2 \cdot (\sigma_t^2(z) p_t(z))$

• Reverse SDE evolution of z_t , $t: 1 \rightarrow 0$ that preserves the density path

 $dz_t = [f_t(z_t) - \sigma_t^2(z_t)\nabla_z \log p_t(z_t)]dt + \sigma_t(z_t)d\overline{W}_t, \quad z_1 \sim p_1(z_1)$

 $\approx \tilde{f}_t(z_t) \text{ ("backward drift")} \qquad \Rightarrow \ f_t(z_t) = \tilde{f}_t(z_t) + \sigma_t^2(z_t) \nabla_z \log p_t(z_t)$

• Re-define the time index: $t \to t/T$, $z_t \to z_{t/T}$, and take limit $T \to \infty$

approx. posterior: $q(z_{[0,1]}|x): q(z_0|x) = N(m, \Sigma)$,

Objective:

• VI-style: $KL[q(z_{[0,1]}|x)||r(z_{[0,1]})]$

• Conditions for a well-defined KL divergence between SDEs:

 $dz_t = \mu_t(z_t, x)dt + \sigma_t d\overline{W}_t, t: 0 \to 1$

- Going in the same direction for analytic form of the KL: reverse either the *r* SDE or the *q* SDE Reverse *r* SDE: $dz_t = [f_t(z_t) + \gamma_t^2 \nabla_z \log r_t(z_t)]dt + \gamma_t d\overline{W}_t, t: 0 \to 1$

He et al. No Trick, No Treat: Pursuits and Challenges Towards Simulation-free Training of Neural Samplers. arXiv:2502.06685 Vargas et al. Denoising Diffusion Samplers. ICLR 2023 Li et al. Scalable Gradients for Stochastic Differential Equations. AISTATS 2020 Song et al. Score-Based Generative Modeling through Stochastic Differential Equations. ICLR 2021

• Re-define the time index: $t \to t/T$, $z_t \to z_{t/T}$, and take limit $T \to \infty$

approx. posterior: $q(z_{[0,1]}|x): q(z_0|x) = N(m, \Sigma)$,

Objective:

- VI-style: $KL[q(z_{[0,1]}|x)||r(z_{[0,1]})]$
- $dz_t = \mu_t(z_t, x)dt + \sigma_t d\overline{W}_t, t: 0 \to 1$
- Conditions for a well-defined KL divergence between SDEs:
 - Going in the same direction for analytic form of the KL: reverse either the *r* SDE or the *q* SDE Reverse *r* SDE: $dz_t = [f_t(z_t) + \gamma_t^2 \nabla_z \log r_t(z_t)]dt + \gamma_t d\overline{W}_t, t: 0 \to 1$
 - Same amount of diffusion: $\sigma_t = \gamma_t$

$$KL[q(z_{[0,1]}|x)||r(z_{[0,1]})] = KL[q(z_0|x)||r(z_0)] + \frac{1}{2}E_q[\int_0^1 \frac{1}{\sigma_t^2} ||f_t(z_t) + \gamma_t^2 \nabla_z \log r_t(z_t) - \mu_t(z_t, x)||_2^2 dt]$$

He et al. No Trick, No Treat: Pursuits and Challenges Towards Simulation-free Training of Neural Samplers. arXiv:2502.06685 Vargas et al. Denoising Diffusion Samplers. ICLR 2023 Li et al. Scalable Gradients for Stochastic Differential Equations. AISTATS 2020 Song et al. Score-Based Generative Modeling through Stochastic Differential Equations. ICLR 2021

• Re-define the time index: $t \to t/T$, $z_t \to z_{t/T}$, and take limit $T \to \infty$

 $\begin{array}{ll} \underline{\text{Diffusion SDE}} & r(z_{[0,1]}) \text{ satisfying } r(z_1) = p(z_1|x), \\ \underline{\text{Augmented}} \text{ model:} & dz_t = \frac{\beta_t z_t dt}{z_t + \sigma \sqrt{2\beta_t} dW_t}, t: 1 \to 0 \\ & (z_1 \coloneqq z) \end{array}$

Objective:

• VI-style: $KL[q(z_{[0,1]}|x)||r(z_{[0,1]})]$

approx. posterior:
$$q(z_{[0,1]}|x): q(z_0|x) = N(m, \Sigma),$$

 $dz_t = \mu_t(z_t, x)dt + \sigma\sqrt{2\beta_t}d\overline{W}_t, t: 0 \to 1$

- Conditions for a well-defined KL divergence between SDEs:
 - Going in the same direction for analytic form of the KL: reverse either the *r* SDE or the *q* SDE Reverse *r* SDE: $dz_t = [\beta_t z_t + 2\sigma^2 \beta_t \nabla_z \log r_t(z_t)]dt + \sigma \sqrt{2\beta_t} d\overline{W}_t, t: 0 \rightarrow 1$
 - Same amount of diffusion: $\sigma_t = \gamma_t = \sigma \sqrt{2\beta_t}$

$$KL[q(z_{[0,1]}|x)||r(z_{[0,1]})] = KL[q(z_0|x)||r(z_0)] + \frac{1}{2\sigma^2}E_q[\int_0^1 \frac{1}{2\beta_t}||\beta_t z_t + 2\sigma^2\beta_t \nabla_z \log r_t(z_t) - \mu_t(z_t, x)||_2^2 dt]$$

He et al. No Trick, No Treat: Pursuits and Challenges Towards Simulation-free Training of Neural Samplers. arXiv:2502.06685

Vargas et al. Denoising Diffusion Samplers. ICLR 2023

Li et al. Scalable Gradients for Stochastic Differential Equations. AISTATS 2020

Song et al. Score-Based Generative Modeling through Stochastic Differential Equations. ICLR 2021

• Re-define the time index: $t \to t/T$, $z_t \to z_{t/T}$, and take limit $T \to \infty$

 $\begin{array}{ll} \underline{\text{Diffusion SDE}} & r(z_{[0,1]}) \text{ satisfying } r(z_1) = p(z_1|x), \\ \underline{\text{Augmented}} \text{ model:} & dz_t = \beta_t z_t dt + \sigma \sqrt{2\beta_t} dW_t, t: 1 \to 0 \\ & (z_1 \coloneqq z) \end{array}$

Objective:

• VI-style: $KL[q(z_{[0,1]}|x)||r(z_{[0,1]})]$

 $\begin{array}{l} \underline{\text{Score param.}}\\ \text{approx. posterior:} \end{array} \begin{array}{l} q\big(z_{[0,1]}|x\big):q(z_0|x)=N(m,\Sigma),\\ dz_t=[\beta_t z_t+2\sigma^2\beta_t s_{\phi}(z_t,x,t)]dt+\sigma\sqrt{2\beta_t}d\overline{W}_t,t:0\to 1 \end{array} \end{array}$

- Conditions for a well-defined KL divergence between SDEs:
 - Going in the same direction for analytic form of the KL: reverse either the *r* SDE or the *q* SDE Reverse *r* SDE: $dz_t = [\beta_t z_t + 2\sigma^2 \beta_t \nabla_z \log r_t(z_t)]dt + \sigma \sqrt{2\beta_t} d\overline{W}_t, t: 0 \to 1$
 - Same amount of diffusion: $\sigma_t = \gamma_t = \sigma \sqrt{2\beta_t}$

$$KL[q(z_{[0,1]}|x)||r(z_{[0,1]})] = KL[q(z_0|x)||r(z_0)] + \sigma^2 E_q[\int_0^1 \beta_t ||\nabla_z \log r_t(z_t) - s_\phi(z_t, x, t)||_2^2 dt]$$

He et al. No Trick, No Treat: Pursuits and Challenges Towards Simulation-free Training of Neural Samplers. arXiv:2502.06685 Vargas et al. Denoising Diffusion Samplers. ICLR 2023 Li et al. Scalable Gradients for Stochastic Differential Equations. AISTATS 2020 Song et al. Score-Based Generative Modeling through Stochastic Differential Equations. ICLR 2021 $\begin{array}{l} \textbf{Estimating the Score} \\ (\text{we can sample } z_{[0,1]} \sim q) \\ KL[q(z_{[0,1]}|x) \| r(z_{[0,1]})] &= KL[q(z_0|x) \| r(z_0)] + \sigma^2 E_q[\int_0^1 \beta_t \| \nabla_z \log r_t(z_t) - s_\phi(z_t, x, t) \|_2^2 dt] \\ \hline \textbf{Diffusion SDE} \\ \textbf{Augmented model:} \\ \begin{array}{l} r(z_{[0,1]}) \text{ satisfying } r(z_1) &= p(z_1|x), \\ dz_t &= \beta_t z_t dt + \sigma \sqrt{2\beta_t} dW_t, t: 1 \rightarrow 0 \\ &= r_t(z_t|z_1) \\ \hline \textbf{Time marginal:} \\ (ignore x notation) \end{array} \\ \begin{array}{l} r_t(z_t) &= \int N(z_t|\sqrt{1 - \lambda_t} z_1, \sigma^2 \lambda_t I)r(z_1) dz_1, \quad \lambda_t = 1 - \exp[-2\int_t^1 \beta_t d\tau] \end{array}$

Estimating the Score (we can sample $z_{[0,1]} \sim q$) $KL[q(z_{[0,1]}|x)||r(z_{[0,1]})] = KL[q(z_0|x)||r(z_0)] + \sigma^2 E_q[\int_0^1 \beta_t ||\nabla_z \log r_t(z_t) - s_\phi(z_t, x, t)||_2^2 dt]$ $r(z_{[0,1]})$ satisfying $r(z_1) = p(z_1|x)$, **Diffusion SDE** <u>Augmented</u> model: $dz_t = \beta_t z_t dt + \sigma_{\sqrt{2\beta_t}} dW_t, t: 1 \to 0$ $\coloneqq r_t(Z_t|Z_1)$ $r_t(z_t) = \int N(z_t) \sqrt{1 - \lambda_t} z_1, \sigma^2 \lambda_t I) r(z_1) dz_1, \quad \lambda_t = 1 - \exp[-2 \int_t^1 \beta_\tau d\tau]$ Time marginal: (ignore x notation) **Denoising Score Identity:** (log derivative trick + Bayes' Rule) $\nabla_{z_t} \log r_t(z_t) = \frac{\int \nabla_{z_t} r_t(z_t|z_1) r(z_1) dz_1}{r_t(z_t)} = \int \frac{r_t(z_t|z_1) r(z_1)}{r_t(z_t)} \nabla_{z_t} \log r_t(z_t|z_1) dz_1 = E_{r_t(z_1|z_t)} [\nabla_{z_t} \log r_t(z_t|z_1)]$

Estimating the Score (we can sample $z_{[0,1]} \sim q$) $KL[q(z_{[0,1]}|x)||r(z_{[0,1]})] = KL[q(z_0|x)||r(z_0)] + \sigma^2 E_q[\int_{1}^{1} \beta_t ||\nabla_z \log r_t(z_t) - s_{\phi}(z_t, x, t)||_2^2 dt]$ $r(z_{[0,1]})$ satisfying $r(z_1) = p(z_1|x)$, **Diffusion SDE** Augmented model: $dz_t = \beta_t z_t dt + \sigma_1 \sqrt{2\beta_t} dW_t, t: 1 \to 0$ $\coloneqq r_t(z_t|z_1)$ $r_t(z_t) = \int N(z_t | \sqrt{1 - \lambda_t} z_1, \sigma^2 \lambda_t I) r(z_1) dz_1, \quad \lambda_t = 1 - \exp[-2 \int_t^1 \beta_\tau d\tau]$ Time marginal: (ignore x notation) Denoising Score Identity: (log derivative trick + Bayes' Rule) $\nabla_{z_t} \log r_t(z_t) = \frac{\int \nabla_{z_t} r_t(z_t|z_1) r(z_1) dz_1}{r_t(z_t)} = \int \frac{r_t(z_t|z_1) r(z_1)}{r_t(z_t)} \nabla_{z_t} \log r_t(z_t|z_1) dz_1 = E_{r_t(z_1|z_t)} [\nabla_{z_t} \log r_t(z_t|z_1)]$ Target Score Identity: $\nabla_{z_t} \log r_t(z_t|z_1) = -\nabla_{z_1} \log r_t(z_t|z_1)$ (Using properties of $r(z_t|z_1)$ as Gaussian) $= -\nabla_{z_1} \log \frac{r_t(z_1|z_t)r_t(z_t)}{r_t(z_1)} = \nabla_{z_1} \log r_1(z_1) - \nabla_{z_1} \log r_t(z_1|z_t)$

He et al. No Trick, No Treat: Pursuits and Challenges Towards Simulation-free Training of Neural Samplers. arXiv:2502.06685

Vargas et al. Denoising Diffusion Samplers. ICLR 2023

De Bortoli et al. Target Score Matching. arXiv:2402.08667

Akhound-Sadegh et al. Iterated Denoising Energy Matching for Sampling from Boltzmann Densities. ICML 2024

Estimating the Score (we can sample $z_{[0,1]} \sim q$) $KL[q(z_{[0,1]}|x)||r(z_{[0,1]})] = KL[q(z_0|x)||r(z_0)] + \sigma^2 E_q[\int_0^1 \beta_t ||\nabla_z \log r_t(z_t) - s_\phi(z_t, x, t)||_2^2 dt]$ $r(z_{[0,1]})$ satisfying $r(z_1) = p(z_1|x)$, **Diffusion SDE** <u>Augmented</u> model: $dz_t = \beta_t z_t dt + \sigma_{\sqrt{2\beta_t}} dW_t, t: 1 \to 0$ $\coloneqq r_t(z_t|z_1)$ $r_t(z_t) = \int N(z_t) \sqrt{1 - \lambda_t} z_1, \sigma^2 \lambda_t I) r(z_1) dz_1, \quad \lambda_t = 1 - \exp[-2 \int_t^1 \beta_\tau d\tau]$ Time marginal: (ignore x notation) Denoising Score Identity: (log derivative trick + Bayes' Rule) $\nabla_{z_t} \log r_t(z_t) = \frac{\int \nabla_{z_t} r_t(z_t|z_1) r(z_1) dz_1}{r_t(z_t)} = \int \frac{r_t(z_t|z_1) r(z_1)}{r_t(z_t)} \nabla_{z_t} \log r_t(z_t|z_1) dz_1 = E_{r_t(z_1|z_t)} [\nabla_{z_t} \log r_t(z_t|z_1)]$ Target Score Identity: $\nabla_{z_t} \log r_t(z_t) = E_{r_t(z_1|z_t,x)} \left[\nabla_{z_1} \log r_1(z_1) - \nabla_{z_1} \log r_t(z_1|z_t) \right] = E_{r_t(z_1|z_t,x)} \left[\nabla_{z_1} \log r_1(z_1) \right]$

He et al. No Trick, No Treat: Pursuits and Challenges Towards Simulation-free Training of Neural Samplers. arXiv:2502.06685 Vargas et al. Denoising Diffusion Samplers. ICLR 2023 De Bortoli et al. Target Score Matching. arXiv:2402.08667 Akhound-Sadegh et al. Iterated Denoising Energy Matching for Sampling from Boltzmann Densities. ICML 2024

Estimating the Score (we can sample $z_{[0,1]} \sim q$) $KL[q(z_{[0,1]}|x)||r(z_{[0,1]})] = KL[q(z_0|x)||r(z_0)] + \sigma^2 E_q[\int_0^1 \beta_t ||\nabla_z \log r_t(z_t) - s_\phi(z_t, x, t)||_2^2 dt]$ $r(z_{[0,1]})$ satisfying $r(z_1) = p(z_1|x)$, **Diffusion SDE** Augmented model: $dz_t = \beta_t z_t dt + \sigma_{\sqrt{2\beta_t}} dW_t, t: 1 \to 0$ $\coloneqq r_t(Z_t|Z_1)$ $r_t(z_t) = \int N(z_t | \sqrt{1 - \lambda_t} z_1, \sigma^2 \lambda_t I) r(z_1) dz_1, \quad \lambda_t = 1 - \exp[-2 \int_t^1 \beta_\tau d\tau]$ Time marginal: (ignore x notation) $\nabla_{z_t} \log r_t(z_t) = E_{r_t(z_1|z_t)} [\nabla_{z_t} \log r_t(z_t|z_1)] = E_{r_t(z_1|z_t)} [\nabla_{z_1} \log r_1(z_1)]$ Denoising Score Identity Target Score Identity

He et al. No Trick, No Treat: Pursuits and Challenges Towards Simulation-free Training of Neural Samplers. arXiv:2502.06685 Vargas et al. Denoising Diffusion Samplers. ICLR 2023 De Bortoli et al. Target Score Matching. arXiv:2402.08667 Akhound-Sadegh et al. Iterated Denoising Energy Matching for Sampling from Boltzmann Densities. ICML 2024

Estimating the Score (we can sample $z_{[0,1]} \sim q$) $KL[q(z_{[0,1]}|x)||r(z_{[0,1]})] = KL[q(z_0|x)||r(z_0)] + \sigma^2 E_q[\int_0^1 \beta_t ||\nabla_z \log r_t(z_t) - s_\phi(z_t, x, t)||_2^2 dt]$ $r(z_{[0,1]})$ satisfying $r(z_1) = p(z_1|x)$, **Diffusion SDE** <u>Augmented</u> model: $dz_t = \beta_t z_t dt + \sigma_{\sqrt{2\beta_t}} dW_t, t: 1 \to 0$ $\coloneqq r_t(Z_t|Z_1)$ $r_t(z_t) = \int N(z_t | \sqrt{1 - \lambda_t} z_1, \sigma^2 \lambda_t I) r(z_1) dz_1, \quad \lambda_t = 1 - \exp[-2 \int_t^1 \beta_\tau d\tau]$ Time marginal: (ignore x notation) $\nabla_{z_t} \log r_t(z_t) = E_{r_t(z_1|z_t)} [\nabla_{z_t} \log r_t(z_t|z_1)] = E_{r_t(z_1|z_t)} [\nabla_{z_1} \log r_1(z_1)]$ Denoising Score Identity Target Score Identity Sampling $z_1 \sim r_t(z_1|z_t)$ via importance sampling: $\propto N\left(z_1; \frac{z_t}{\sqrt{1-\lambda_t}}, \frac{\sigma^2 \lambda_t}{1-\lambda_t}I\right) =: n_t(z_1|z_t)$ $\coloneqq r(z_t|z_1)$ $r_t(z_1|z_t) \propto r_1(z_1) \exp\left[-\frac{1}{2\sigma^2 \lambda_t} \left\|z_t - \sqrt{1 - \lambda_t} z_1\right\|_2^2\right] \propto r_1(z_1) \exp\left[-\frac{1 - \lambda_t}{2\sigma^2 \lambda_t} \left\|z_1 - \frac{z_t}{\sqrt{1 - \lambda_t}}\right\|_2^2\right]$

He et al. No Trick, No Treat: Pursuits and Challenges Towards Simulation-free Training of Neural Samplers. arXiv:2502.06685

Vargas et al. Denoising Diffusion Samplers. ICLR 2023

De Bortoli et al. Target Score Matching. arXiv:2402.08667

Akhound-Sadegh et al. Iterated Denoising Energy Matching for Sampling from Boltzmann Densities. ICML 2024

Estimating the Score (we can sample $z_{[0,1]} \sim q$) $KL[q(z_{[0,1]}|x)||r(z_{[0,1]})] = KL[q(z_0|x)||r(z_0)] + \sigma^2 E_q[\int_{1}^{1} \beta_t ||\nabla_z \log r_t(z_t) - s_{\phi}(z_t, x, t)||_2^2 dt]$ $r(z_{[0,1]})$ satisfying $r(z_1) = p(z_1|x)$, **Diffusion SDE** Augmented model: $dz_t = \beta_t z_t dt + \sigma_1 \sqrt{2\beta_t} dW_t, t: 1 \to 0$ $\coloneqq r_t(Z_t|Z_1)$ $r_t(z_t) = \int N(z_t) \sqrt{1 - \lambda_t} z_1, \sigma^2 \lambda_t I) r(z_1) dz_1, \quad \lambda_t = 1 - \exp[-2 \int_t^1 \beta_\tau d\tau]$ Time marginal: (ignore x notation) $\nabla_{z_t} \log r_t(z_t) = E_{r_t(z_1|z_t)} [\nabla_{z_t} \log r_t(z_t|z_1)] = E_{r_t(z_1|z_t)} |\nabla_{z_1} \log r_1(z_1)|$ Denoising Score Identity Target Score Identity Sampling $z_1 \sim r_t(z_1|z_t)$ via importance sampling: $N\left(z_1; \frac{z_t}{\sqrt{1-\lambda_t}}, \frac{\sigma^2 \lambda_t}{1-\lambda_t} I\right) =: n_t(z_1|z_t)$ $\nabla_{z_t} \log r_t(z_t) = E_{n_t(z_1|z_t)}[w(z_1)\nabla_{z_t} \log r_t(z_t|z_1)] = E_{n_t(z_1|z_t)}[w(z_1)\nabla_{z_1} \log r_1(z_1)]$ $w(z_1) \propto \tilde{r}_1(z_1)$ (unnormalized target density) He et al. No Trick, No Treat: Pursuits and Challenges Towards Simulation-free Training of Neural Samplers. arXiv:2502.06685

Vargas et al. Denoising Diffusion Samplers. ICLR 2023

De Bortoli et al. Target Score Matching. arXiv:2402.08667

Akhound-Sadegh et al. Iterated Denoising Energy Matching for Sampling from Boltzmann Densities. ICML 2024

Auxiliary VI → Diffusion Posterior (Continuous Time)

A summary of the Diffusion algorithm applied to VI

 $\begin{array}{ll} \underline{\text{Diffusion SDE}} & r(z_{[0,1]}) \text{ satisfying } r(z_1) = p(z_1|x) \text{ (target density)} \\ \underline{\text{Augmented model:}} & dz_t = \beta_t z_t dt + \sigma \sqrt{2\beta_t} dW_t, t: 1 \rightarrow 0 \\ & (z_1 \coloneqq z) \end{array}$

 $\begin{array}{ll} \underline{\text{Score param.}} & q\big(z_{[0,1]}|x\big):q(z_0|x)=N(m,\Sigma),\\ \text{approx. posterior:} & dz_t=[\beta_t z_t+2\sigma^2\beta_t s_\phi(z_t,t)]dt+\sigma\sqrt{2\beta_t}d\overline{W}_t,t:0\to 1 \end{array} \end{array}$

• Objective: VI style, $KL[q(z_{[0,1]}|x)||r(z_{[0,1]})]$:

Fit $q(z_0|x)$, $s_{\phi}(z_t, x, t)$ by min. $KL[q(z_0|x)||r(z_0)]$ plus $E_q[\int_0^1 \beta_t ||\nabla_z \log r_t(z_t)| - s_{\phi}(z_t, x, t)||_2^2 dt]$ (estimated via importance sampling + denoising/target score identity) • After fitting the posterior (i.e., $s_{\phi}(z_t, t)$):

Sample $z_0 \sim q(z_0|x)$, and simulate the approx. posterior SDE to get approximately $z_1 \sim p(z_1|x)$

He et al. No Trick, No Treat: Pursuits and Challenges Towards Simulation-free Training of Neural Samplers. arXiv:2502.06685 Vargas et al. Denoising Diffusion Samplers. ICLR 2023 Akhound-Sadegh et al. Iterated Denoising Energy Matching for Sampling from Boltzmann Densities. ICML 2024



Example: Molecular Dynamics Simulations

Benchmark: Lennard-Jones Potential (LJ-N)

- Goal: study intermolecular interactions
- A molecule is a system of *N* atoms ("particles") $z = [z_1, ..., z_N]$
- Each particle z_i has its own 3D coordinates (so dim(z) = 3N) $z_i = [z_i^x, z_i^y, z_i^z]$
- Define $r_{ij} = ||z_i z_j||_2^2$, the Lennard-Jones potential is:

$$V_{LJ}(r_{ij}, n, m) \coloneqq \frac{A_n}{r_{ij}^n} - \frac{B_m}{r_{ij}^m}$$



$$V_{LJ}(r_{ij}) \coloneqq 4\epsilon \left[\frac{\sigma^{12}}{r_{ij}^{12}} - \frac{\sigma^6}{r_{ij}^6} \right] \qquad (\min_r V_{LJ}(r) = -\epsilon)$$

Target distribution: $\pi(z) \propto \exp[-\frac{E(z)}{k_BT}]$, $E(z) \coloneqq \sum_{i,j \neq i} V_{LJ}(r_{ij})$




Example: Molecular Dynamics Simulations

LJ13: Lennard-Jones 12-6 potential, N = 13 particles (dim(z) = 39)



Example: Molecular Dynamics Simulations

Sampling conformations based on energy functions in (quantum) density functional theory (DFT)

Conformers of a molecule:

- spatial arrangements determined by locally stable configurations of rotations around its bonds
 - i.e., local minima on the molecule's potential energy surface
- can be interconverted by rotations around single bonds





Summary (Diffusion-based Approximations)

Constructing flexible q distribution via mixtures:

- Hierarchical mixtures (very flexible)
- The VI objective needs to be augmented (adding in the r distributions)
- Add many layers and take continuous-limit: Diffusion-based approximations
 - Be careful about the validity of the KL divergence definition
 - Similar to diffusion generative model training except for the score estimation method



Designing q Distributions





Normalizing flows

Auxiliary variables & mixture distributions

Use many layers (hierarchical) + continuous-time limit

Diffusion SDE & Continuous Normalizing Flow posteriors for approximate inference

Normalizing Flows (Discrete Time)

- Change-of-variable formula:
 - x is a random variable with probability density function (PDF) $p_X(x)$
 - y = f(x) is an invertible mapping
 - The probability mass is preserved, and the PDF for y = f(x) satisfies

$$p_Y(y)dy = p_X(x)dx$$

prob. mass of region around y



$$p_Y(y) = p_X(x) |\det(\frac{dx}{dy})|$$
$$p_X(x) = p_Y(y) |\det(\frac{dy}{dx})|$$



Normalizing Flows (Discrete Time)

- Variational inference with normalizing flow
 - Assume $q_0(z_0) = N(z_0; 0, I)$
 - Define $z = T_{\phi}(z_0)$ where $T_{\phi}(\cdot)$ is an invertible mapping parameterized by ϕ

$$q(z) = q_0(z_0) |\det\left(\frac{dz}{dz_0}\right)|^{-1}$$
 with $z_0 = T_{\phi}^{-1}(z)$

• Fit q(z) to p(x | z) with VI:

$$E_{q(z)} = E_{q(z)} [\log p(x \mid z) + \log p(z) - \log q(z)]$$
 by def. of $q(z)$
= $E_{q(z)} \left[\log p(x, z) - \left[\log q_0(z_0 = T_{\phi}^{-1}(z)) \right] \det \left(\frac{dz}{dz_0} \right) \right]^{-1} = E_{q_0(z_0)} \left[\log p(x, T_{\phi}(z_0)) - \log q_0(z_0) + \log \left| \det \left(\frac{df_{\phi}}{dz_0} \right) \right| \right]$

reparam. trick:
$$z \sim q(z) \Leftrightarrow z_0 \sim q_0(z_0), z = T_{\phi}(z_0)$$

• Computing ELBO requires $\log |\det \left(\frac{dT_{\phi}}{dz_0}\right)|$

Rezende and Mohamed. Variational Inference with Normalizing Flows. ICML 2015

Normalizing Flows (Discrete Time)

- Variational inference with normalizing flow
 - Idea: define T_{ϕ} such that $\log |\det \left(\frac{dT_{\phi}}{dz_0}\right)|$ is easy to compute!
 - Chain simple invertible mappings together to make a flexible mapping



$$T_{\phi} = f_K \circ f_{K-1} \circ \cdots \circ f_1, f_k(\cdot) \coloneqq f_{\phi_k}(\cdot), \phi = \{\phi_k\}_{k=1}^K$$

• For each mapping, hopefully the Jacobian log-determinant is easy to compute

$$\Rightarrow \log |\det\left(\frac{dT_{\phi}}{dz_{0}}\right)| = \sum_{k=1}^{K} \log |\det\left(\frac{dz_{k}}{dz_{k-1}}\right)|$$

How about the continuous-time limit? $(k \rightarrow k/K, K \rightarrow \infty)$

Rezende and Mohamed. Variational Inference with Normalizing Flows. ICML 2015

Continuous-Time Limit: Continuous Normalizing Flows

• Continuous Normalizing Flows (CNFs):

$$z_1 = T(z_0), T(z_0) \coloneqq z_0 + \int_0^1 v_t(z_t) dt$$

- Motivation: Enhanced expressiveness
- Change-of-variable rule:

$$\partial_t \log p_t(z_t) = -\nabla_z \cdot v_t(z_t)$$
$$\log p_1(z_1) = \log p_0(z_0) - \int_0^1 \nabla_z \cdot v_t(z_t) dt$$

(expensive to compute!)



Continuity Equation

• Forward SDE evolution of z_t , $t: 0 \rightarrow 1$

 $dz_t = v_t(z_t)dt + \sigma_t(z_t)dW_t, \quad z_0 \sim p_0(z_0)$

• Fokker-Planck Equation: marginal density evolution of $p_t(z), t: 0 \to 1$, for any z (not necessarily for $z = z_t$) $\partial_t p_t(z) = -\nabla_z \cdot \left(v_t(z) p_t(z) \right) + \frac{1}{2} \nabla_z^2 \cdot \left(\sigma_t^2(z) p_t(z) \right)$

Continuity Equation (Cont.)

• Forward ODE evolution of z_t , $t: 0 \rightarrow 1$

$$dz_t = v_t(z_t)dt + \sigma_t(z_t)dW_t, \quad z_0 \sim p_0(z_0)$$

• Fokker-Planck Equation: marginal density evolution of $p_t(z), t: 0 \to 1$, for any z (not necessarily for $z = z_t$) $\partial_t p_t(z) = -\nabla_z \cdot \left(v_t(z)p_t(z)\right) + \frac{1}{2}\nabla_z^2 \cdot (o_t^2(z)p_t(z))$

Divide $p_t(z)$ on both sides: $\partial_t \log p_t(z) = -\nabla_z \cdot v_t(z) - \langle \nabla_z \log p_t(z), v_t(z) \rangle$

• Straight-forward reverse simulation of z_t , $t: 1 \rightarrow 0$

 $dz_t = v_t(z_t)dt, \quad z_1 \sim p_1(z_1) \qquad \Rightarrow \quad z_{t-\delta_t} \approx z_t - v_t(z_t)\delta t$

Continuous-Time Limit: Continuous Normalizing Flows

• Continuous Normalizing Flows (CNFs):

$$z_1 = T(z_0), T(z_0) \coloneqq z_0 + \int_0^1 v_t(z_t) dt$$

• Probability density evolves: $\{p_t(z)\}_{t \in [0,1]}$ satisfy

$$\partial_t \log p_t(\mathbf{z}) = -\nabla_{\mathbf{z}} \cdot v_t(\mathbf{z}) - \langle \nabla_{\mathbf{z}} \log p_t(\mathbf{z}), v_t(\mathbf{z}) \rangle$$

• Note the difference from change-of-variable rule:

$$\partial_t \log p_t(z_t) = -\nabla_z \cdot v_t(z_t)$$



• Specify a density path:

 $\{p_t(z)\}_{t\in[0,1]}$: $p_0(z)$ easy to sample,

$$p_t(z) = \frac{1}{z_t} \exp[-E_t(z)],$$

$$p_1(z) = \pi(z), \text{ i.e., } E_1(z) \coloneqq E(z)$$

E.g., tempering:

 $E_t(z) = \beta_t E_0(z) + (1 - \beta_t) E(z),$

 $\pi(z) = \frac{1}{7} \exp[-E(z)]$

 $\beta_0 = 1, \beta_1 = 0$

• Specify a density path:

 $\{p_t(z)\}_{t\in[0,1]}: \qquad p_t(z) = \frac{1}{Z_t} \exp[-E_t(z)],$ $p_0(z) \text{ easy to sample,} \qquad p_1(z) = \pi(z), \text{ i.e., } E_1(z) \coloneqq E(z)$

• Learn a CNF network $v_{\phi}(z, t)$ by minimizing L2 error ("PINN loss"):

 $L(v_{\phi}) \coloneqq E_{q_t(z)}[\|\partial_t \log p_t(z) + \nabla_z \cdot v_{\phi}(z, t) + \langle \nabla_z \log p_t(z), v_{\phi}(z, t) \rangle \|_2^2]$

Ensuring continuity equation to hold for every $z \sim q_t(z)$

E.g., tempering:

 $E_t(z) = \beta_t E_0(z) + (1 - \beta_t) E(z),$

Tian et al. Liouville Flow Importance Sampler. ICML 2024 Mate and Fleuret. Learning Interpolations between Boltzmann Densities. TMLR 2023 Chen et al. Neural Flow Samplers with Shortcut Models. arXiv:2502.07337 $\pi(z) = \frac{1}{7} \exp[-E(z)]$

 $\beta_0 = 1, \beta_1 = 0$

• Specify a density path:

 $\{p_t(z)\}_{t\in[0,1]}: \qquad p_t(z) = \frac{1}{Z_t} \exp[-E_t(z)],$ $p_0(z) \text{ easy to sample,} \qquad p_1(z) = \pi(z), \text{ i.e., } E_1(z) \coloneqq E(z)$

• Learn a CNF network $v_{\phi}(z, t)$ by minimizing L2 error ("PINN loss"):

 $L(v_{\phi}) \coloneqq E_{q_t(z)}[\| \partial_t \log p_t(z) + \nabla_z \cdot v_{\phi}(z, t) + \langle \nabla_z \log p_t(z), v_{\phi}(z, t) \rangle \|_2^2]$

Ensuring continuity equation to hold for every $z \sim q_t(z)$

E.g., tempering:

 $E_t(z) = \beta_t E_0(z) + (1 - \beta_t) E(z),$

• Simulate samples from $\pi(z)$ (approximately) by solving ODE:

$$x_0 \sim p_0(z), \qquad z_1 \coloneqq z_0 + \int_0^1 v_{\phi}(z_t, t) dt$$

Tian et al. Liouville Flow Importance Sampler. ICML 2024

Mate and Fleuret. Learning Interpolations between Boltzmann Densities. TMLR 2023 Chen et al. Neural Flow Samplers with Shortcut Models. arXiv:2502.07337 $\pi(z) = \frac{1}{7} \exp[-E(z)]$

 $\beta_0 = 1, \beta_1 = 0$

$\pi(z) = \frac{1}{z} \exp[-E(z)]$ $= -\nabla_z E_t(z)$ Training: $L(v_{\phi}) \coloneqq E_{q_t(z)}[\| \partial_t \log p_t(z) + \nabla_z \cdot v_{\phi}(z, t) + \langle \nabla_z \log p_t(z) \rangle, v_{\phi}(z, t) \rangle \|_2^2]$ Sampling: $z_0 \sim p_0(z), \quad z_1 \coloneqq z_0 + \int_0^1 v_{\phi}(z_t, t) dt$

- Challenges:
 - Selecting the "training distribution" $q_t(z)$ and estimating the expectation
 - Not necessary for $q_t(z) = p_t(z)$ but ideally $q_t(z) \approx p_t(z)$

$\pi(z) = \frac{1}{z} \exp[-E(z)]$ $= -\nabla_z E_t(z)$ Training: $L(v_{\phi}) \coloneqq E_{q_t(z)}[\| \partial_t \log p_t(z) + \nabla_z \cdot v_{\phi}(z, t) + \langle \nabla_z \log p_t(z) \rangle, v_{\phi}(z, t) \rangle \|_2^2]$ Sampling: $z_0 \sim p_0(z), \quad z_1 \coloneqq z_0 + \int_0^1 v_{\phi}(z_t, t) dt$

- Challenges:
 - Selecting the "training distribution" $q_t(z)$ and estimating the expectation
 - Not necessary for $q_t(z) = p_t(z)$ but ideally $q_t(z) \approx p_t(z)$
 - Estimating $\partial_t \log p_t(z)$:

$$p_t(z) \coloneqq \frac{1}{Z_t} \exp[-E_t(z)] \quad \Rightarrow \quad \partial_t \log p_t(z) = -\partial_t E_t(z) - \frac{\partial_t \log Z_t}{(\text{intractable})}$$

Chen et al. Neural Flow Samplers with Shortcut Models. arXiv:2502.07337

$\pi(z) = \frac{1}{Z} \exp[-E(z)]$ $= -\nabla_z E_t(z)$ Training: $L(v_{\phi}) \coloneqq E_{q_t(z)}[\| \partial_t \log p_t(z) + \nabla_z \cdot v_{\phi}(z, t) + \langle \nabla_z \log p_t(z) \rangle, v_{\phi}(z, t) \rangle \|_2^2]$ Sampling: $z_0 \sim p_0(z), \quad z_1 \coloneqq z_0 + \int_0^1 v_{\phi}(z_t, t) dt$

- Challenges:
 - Selecting the "training distribution" $q_t(z)$ and estimating the expectation
 - Not necessary for $q_t(z) = p_t(z)$ but ideally $q_t(z) \approx p_t(z)$
 - Estimating $\partial_t \log p_t(z)$:

$$p_t(z) \coloneqq \frac{1}{Z_t} \exp[-E_t(z)] \implies \partial_t \log p_t(z) = -\partial_t E_t(z) - \partial_t \log Z_t$$

(intractable)

• Solving the ODE flow simulation in a fast way

Using "Training Data" $q_t(z)$

 $L(v_{\phi}) \coloneqq E_{\boldsymbol{q}_{t}(\boldsymbol{z})}[\| \partial_{t} \log p_{t}(\boldsymbol{z}) + \nabla_{\boldsymbol{z}} \cdot v_{\phi}(\boldsymbol{z}, t) + \langle \nabla_{\boldsymbol{z}} \log p_{t}(\boldsymbol{z}), v_{\phi}(\boldsymbol{z}, t) \rangle \|_{2}^{2}]$

- Sampling under $q_t(z) \approx p_t(z)$ via velocity-driven SMC:
 - A typical SMC method (e.g., Hamiltonian AIS):
 - Pick $0 = t_0 < t_1 < t_2 < \dots < t_M = 1$ and run SMC with path $\{p_{t_m}(z)\}_{m=0}^M$ as proposals
 - Compute the importance weights by accumulating density ratios through time
 - Resampling is required by monitoring ESS
 - The steps for approximately drawing samples from $p_{t_m}(z)$:
 - Transport from previous step: $\tilde{z}_{t_m} = z_{t_{m-1}} + \int_{t_{m-1}}^{t_m} v_{\phi}(z_t, t) dt$ ("prediction")
 - Run (short-chain) HMC: $z_{t_m} = HMC(\tilde{z}_{t_m})$ ("correction")

Neal. Annealed Importance Sampling. Stats. Comp., 2001 Neal. MCMC using Hamiltonian Dynamics. Handbook of MCMC, 2010 Chen et al. Neural Flow Samplers with Shortcut Models. arXiv:2502.07337

 $\pi(z) = \frac{1}{Z} \exp[-E(z)]$

Estimating
$$\partial_t \log Z_t$$

Monte Carlo estimate can have high variance!

$$\partial_t \log Z_t = \frac{1}{Z_t} \partial_t \int \exp[-E_t(z)] \, dz = -\frac{1}{Z_t} \int \exp[-E_t(z)] \, \partial_t E_t(z) \, dz = -\frac{E_{p_t(z)}}{Z_t} [\partial_t E_t(z)] \, dz$$

• Solution: Stein control variate with $z^k \sim p_t(z)$

(Langevin-Stein Operator)

 $-E_{p_t(z)}[\partial_t E(z)] \approx \frac{1}{K} \sum_{k=1}^K -\partial_t E_t(z^k) + \beta \left[\nabla_z \cdot v_\phi(z^k, t) + \langle \nabla_x \log p_t(z^k), v_\phi(z^k, t) \rangle \right]$

Variance Reduction for Monte Carlo

- Control variate method:
 - Assume we want to estimate with MC simulation

 $E_{q(z)}[F(z)] \approx \frac{1}{K} \sum_{k=1}^{K} F(z^k), \quad z^k \sim q(z)$

- Control variate: define a control function G(z) satisfying:
 - $V_{q(z)}[G(z)] < \infty$
 - Known or fast computable $E_{q(z)}[G(z)]$



G(z)

Variance Reduction for Monte Carlo

- Control variate method:
 - Then define the new MC estimator

$$E_{q(z)}[F(z)] \approx \frac{1}{K} \sum_{k=1}^{K} \widehat{F}(z^k), \ z^k \sim q(z),$$



< 0 if F and G are strongly and positively correlated

 $\pi(z) = \frac{1}{Z} \exp[-E(z)]$

Estimating
$$\partial_t \log Z_t$$

Monte Carlo estimate can have high variance!

$$\partial_t \log Z_t = \frac{1}{Z_t} \partial_t \int \exp[-E_t(z)] \, dz = -\frac{1}{Z_t} \int \exp[-E_t(z)] \, \partial_t E_t(z) \, dz = -\frac{E_{p_t(z)}}{Z_t} [\partial_t E_t(z)]$$

• Solution: Stein control variate with $z^k \sim p_t(z)$

(Langevin-Stein Operator)

 $-E_{p_t(z)}[\partial_t E(z)] \approx \frac{1}{K} \sum_{k=1}^K -\partial_t E_t(z^k) + \beta \left[\nabla_z \cdot v_\phi(z^k, t) + \langle \nabla_x \log p_t(z^k), v_\phi(z^k, t) \rangle \right]$

- Stein's Identity ensures unbiasedness: for any $v_{\phi}(z, t)$ $E_{p_t(z)} [\nabla_z \cdot v_{\phi}(z, t) + \langle \nabla_z \log p_t(z), v_{\phi}(z, t) \rangle] = 0$
- Continuity equation for globally optimal $v_{\phi}(z,t)$

 $\frac{-\partial_t E_t(z) + \left[\nabla_z \cdot v_\phi(z, t) + \left\langle \nabla_z \log p_t(z), v_\phi(z, t) \right\rangle\right]}{\coloneqq \xi_t(z; v_\phi(z, t))} = \partial_t \log Z_t$

 \Rightarrow Variance is 0 when $\beta = 1$ and $v_{\phi}(z, t)$ is optimal

Liu et al. Action-dependent Control Variates for Policy Optimization via Stein Identity. ICLR 2018 Chen et al. Neural Flow Samplers with Shortcut Models. arXiv:2502.07337



Speeding up with Shortcuts

Sampling: $z_0 \sim p_0(z)$, $z_1 \coloneqq z_0 + \int_0^1 v_{\phi}(z_t, t) dt$

Naïve Euler method requires a lot of steps!

 $\pi(z) = \frac{1}{7} \exp[-E(z)]$

Solution: Shortcut models

$$z_{t+d} \leftarrow z_t + s_{\phi}(z_t, t, d)d$$

- A shortcut model s_{θ} is a valid ODE solver if for any z_t :
 - $s_{\phi}(z_t, t, 0) \coloneqq v_{\phi}(z_t, t)$

Frans et al. One Step Diffusion via Shortcut Models. ICLR 2025

- $s_{\phi}(z_t, t, 2d) = \frac{1}{2}s_{\phi}(z_t, t, d) + \frac{1}{2}s_{\phi}(z_{t+d}, t+d, d)$
- Training a neural flow shortcut sampler:

$$\tilde{L}(s_{\phi}) \coloneqq L\left(s_{\phi}(\cdot, \cdot, 0)\right) + E_{q(z_{t}, d)}\left[\|s_{\phi}(z_{t}, t, 2d) - \frac{1}{2}s_{\phi}(z_{t}, t, d) - \frac{1}{2}s_{\phi}(z_{t+d}, t+d, d)\|_{2}^{2}\right]$$

flow learning enforcing consistency



Target $\pi(z)$: mixture of 40 Gaussians



- "Ground Truth": samples from mixture of Gaussians
- FAB: normalising flow transport map, trained by alpha-divergence, "data" from AIS + replay buffer
- iDEM: diffusion-based, score estimation via importance sampling + replay buffer
- LFIS: continuity equation-based loss, no amortization across t, simple importance sampling for $\partial_t \log Z_t$

Midgley et al. Flow Annealed Importance Sampling Bootstrap. ICLR 2023 Akhound-Sadegh et al. Iterated Denoising Energy Matching for Sampling from Boltzmann Densities. ICML 2024 Tian et al. Liouville Flow Importance Sampler. ICML 2024 Chen et al. Neural Flow Samplers with Shortcut Models. arXiv:2502.07337

Target $\pi(z)$: mixture of 40 Gaussians



Target $\pi(z)$: mixture of 40 Gaussians

Ablation: Training with or without the shortcut consistency loss

$$E(s_{\phi}) \coloneqq L(s_{\phi}(\cdot, \cdot, 0)) + E_{q(z_{t}, d)} \left[\|s_{\phi}(z_{t}, t, 2d) - \frac{1}{2}s_{\phi}(z_{t}, t, d) - \frac{1}{2}s_{\phi}(z_{t+d}, t+d, d)\|_{2}^{2} \right]$$

flow learning

enforcing consistency



Target $\pi(z)$: MW32 (32D Many-Well potential with $2^{16} = 65,536$ modes) (only showing two dimensions here)



Target $\pi(z)$: LJ-13 (39D, energy dependent on atom distances)

$$\pi(z) \propto \exp\left[-\frac{E(z)}{k_BT}\right], E(z) \coloneqq \sum_{i,j\neq i} V_{LJ}(r_{ij}), r_{ij} = \left\|z_i - z_j\right\|_2^2$$
$$V_{LJ}(r_{ij}) \coloneqq 4\epsilon \left[\frac{\sigma^{12}}{r_{ij}^{12}} - \frac{\sigma^6}{r_{ij}^6}\right]$$

Method	Energy W_2	Energy TV	Distance TV
FAB	31.28 ± 0.31	0.94 ± 0.03	0.26 ± 0.01
iDEM	13.13 ± 5.30	0.31 ± 0.01	0.03 ± 0.01
LFIS	∞	*	0.88 ± 0.00
LIBD	49.89 ± 0.12	*	0.45 ± 0.01
PINN	48.3 ± 0.1	*	0.44 ± 0.00
NFS ² -128 (ours)	1.93 ± 0.01	0.19 ± 0.10	0.05 ± 0.00
NFS 2 -64 (ours)	1.51 ± 0.10	0.19 ± 0.10	0.06 ± 0.01
NFS ² -32 (ours)	1.62 ± 0.20	0.20 ± 0.01	0.06 ± 0.00
NFS 2 -8 (ours)	29.57 ± 16.06	0.30 ± 0.01	0.22 ± 0.01

Chen et al. Neural Flow Samplers with Shortcut Models. arXiv:2502.07337



Summary (CNF-based Approximations)

Constructing flexible *q* distribution via transport:

- Normalizing flows as stacking invertible transformations (very flexible)
- The VI objective needs to compute $\log |\det \left(\frac{dT}{dz}\right)|$
- Add many layers and take continuous-limit: CNF-based approximations
 - The log-determinant requires explicitly integrating an ODE (expensive)
 - "simulation-free" approach by leveraging the continuity equation ("PINN loss")
 - Many challenges need solutions

Training:
$$L(v_{\phi}) \coloneqq E_{q_t(z)}[\|\partial_t \log p_t(z) + \nabla_z \cdot v_{\phi}(z,t) + \langle \nabla_z \log p_t(z), v_{\phi}(z,t) \rangle \|_2^2]$$

Sampling:
$$z_0 \sim p_0(z), \qquad z_1 \coloneqq z_0 + \int_0^1 v_\phi(z_t, t) dt$$

Designing q Distributions





Normalizing flows

Auxiliary variables & mixture distributions

Use many layers (hierarchical) + continuous-time limit

Diffusion SDE & Continuous Normalizing Flow posteriors for approximate inference Converting between each other?

• Specify a density path:

 $\{p_t(z)\}_{t \in [0,1]}: \qquad p_t(z) = \frac{1}{Z_t} \exp[-E_t(z)],$ $p_0(z) \text{ easy to sample,} \qquad p_1(z) = \pi(z), \text{ i.e., } E_1(z) \coloneqq E(z)$

• Learn a CNF network $v_{\phi}(z, t)$ by minimizing L2 error ("PINN loss"):

 $L(v_{\phi}) \coloneqq E_{q_t(z)}[\| \partial_t \log p_t(z) + \nabla_x \cdot v_{\phi}(z, t) + \langle \nabla_z \log p_t(z), v_{\phi}(z, t) \rangle \|_2^2]$

Ensuring continuity equation to hold for every $z \sim q_t(z)$

• Simulate samples from $\pi(z)$ (approximately) by solving ODE:

$$x_0 \sim p_0(z), \qquad z_1 \coloneqq z_0 + \int_0^1 v_{\phi}(z_t, t) dt$$

This path can also be defined via diffusion!

 $dz_t = \beta_t z_t dt + \sigma \sqrt{2\beta_t} dW_t, t: 1 \to 0$

• Specify a density path:

 $\{p_t(z)\}_{t\in[0,1]}: \qquad p_t(z) = \frac{1}{Z_t} \exp[-E_t(z)],$ $p_0(z) \text{ easy to sample,} \qquad p_1(z) = \pi(z), \text{ i.e., } E_1(z) \coloneqq E(z)$

• Learn a CNF network $v_{\phi}(z, t)$ by minimizing L2 error ("PINN loss"): require estimations $L(v_{\phi}) \coloneqq E_{q_t(z)}[\|\partial_t \log p_t(z)| + \nabla_x \cdot v_{\phi}(z, t) + \langle \nabla_z \log p_t(z)|, v_{\phi}(z, t) \rangle \|_2^2]$

Ensuring continuity equation to hold for every $z \sim q_t(z)$

• Simulate samples from $\pi(z)$ (approximately) by solving ODE:

$$x_0 \sim p_0(z), \qquad z_1 \coloneqq z_0 + \int_0^1 v_{\phi}(z_t, t) dt$$

This path can also be defined via diffusion!

• Specify a density path: $dz_t = \beta_t z_t dt + \sigma \sqrt{2\beta_t} dW_t, t: 1 \to 0$ ${p_t(z)}_{t \in [0,1]}$: $p_t(z) = \frac{1}{Z_t} \exp[-E_t(z)],$ $p_0(z)$ easy to sample, $p_1(z) = \pi(z)$, i.e., $E_1(z) \coloneqq E(z)$ $\coloneqq p_t(z_t|z_1)$ $p_t(z_t) = \int N(z_t | \sqrt{1 - \lambda_t} z_1, \sigma^2 \lambda_t I) p_1(z_1) dz_1, \quad \lambda_t = 1 - \exp[-2\int_t^1 \beta_\tau d\tau]$ Time marginal: Estimate the score $\nabla_{z_t} \log p_t(z_t)$ via importance sampling: $N\left(z_1; \frac{z_t}{\sqrt{1-\lambda_t}}, \frac{\sigma^2 \lambda_t}{1-\lambda_t}I\right) \coloneqq n_t(z_1|z_t)$ Denoising Score Identity Target Score Identity $\nabla_{z_t} \log p_t(z_t) = E_{n_t(z_1|z_t)}[w(z_1)\nabla_{z_t} \log p_t(z_t|z_1)] = E_{n_t(z_1|z_t)}[w(z_1)\nabla_{z_1} \log r_1(z_1)]$ $w(z_1) \propto \exp[-E_1(z_1)]$ (unnormalized target density)

This path can also be defined via diffusion!

• Specify a density path: $dz_t = \beta_t z_t dt + \sigma_1 \sqrt{2\beta_t} dW_t, t: 1 \to 0$ ${p_t(z)}_{t \in [0,1]}$: $p_t(z) = \frac{1}{Z_t} \exp[-E_t(z)],$ $p_0(z)$ easy to sample, $p_1(z) = \pi(z)$, i.e., $E_1(z) \coloneqq E(z)$ $\coloneqq p_t(z_t|z_1)$ $p_t(z_t) = \int N(z_t | \sqrt{1 - \lambda_t} z_1, \sigma^2 \lambda_t I) p_1(z_1) dz_1, \quad \lambda_t = 1 - \exp[-2\int_t^1 \beta_\tau d\tau]$ Time marginal: Estimate the score $\nabla_{z_t} \log p_t(z_t)$ via importance sampling: $N\left(z_1; \frac{z_t}{\sqrt{1-\lambda_t}}, \frac{\sigma^2 \lambda_t}{1-\lambda_t}I\right) \coloneqq n_t(z_1|z_t)$ Denoising Score Identity Target Score Identity $\nabla_{z_t} \log p_t(z_t) = E_{n_t(z_1|z_t)}[w(z_1)\nabla_{z_t} \log p_t(z_t|z_1)] = E_{n_t(z_1|z_t)}[w(z_1)\nabla_{z_1} \log r_1(z_1)]$ $w(z_1) \propto \exp[-E_1(z_1)]$ (unnormalized target density) Estimate the "time-score" $\partial_t \log p_t(z_t)$ also via importance sampling:

 $\partial_t \log p_t(z_t) = E_{n_t(z_1|z_t)}[w(z_1)\partial_t \log p_t(z_t|z_1)]$

Converting Between ODEs and SDEs

- Forward ODE/SDE evolution of z_t , $t: 0 \rightarrow 1$, $z_0 \sim p_0(z_0)$ CNF-based sampler (ODE) Diffusion-based sampler (SDE) $dz_t = v_t(z_t)dt$ $dz_t = f_t(z_t)dt + \sigma_t dW_t$
- Fokker-Planck Equation: marginal density evolution of $p_t(z), t: 0 \rightarrow 1$, for any z (not necessarily for $z = z_t$)

CNF-based sampler (ODE)Diffusion-based sampler (SDE) $\partial_t p_t(z) = -\nabla_z \cdot \left(v_t(z) p_t(z) \right)$ $\partial_t p_t(z) = -\nabla_z \cdot \left(f_t(z) p_t(z) \right) + \frac{1}{2} \nabla_z^2 \cdot (\sigma_t^2 p_t(z))$

Converting Between ODEs and SDEs

- Forward ODE/SDE evolution of z_t , $t: 0 \rightarrow 1$, $z_0 \sim p_0(z_0)$ CNF-based sampler (ODE) Diffusion-based sampler (SDE) $dz_t = v_t(z_t)dt$ $dz_t = f_t(z_t)dt + \sigma_t dW_t$
- Fokker-Planck Equation: marginal density evolution of $p_t(z), t: 0 \rightarrow 1$, for any z (not necessarily for $z = z_t$)
Converting Between ODEs and SDEs

- Forward ODE/SDE evolution of z_t , $t: 0 \rightarrow 1$, $z_0 \sim p_0(z_0)$ CNF-based sampler (ODE) Diffusion-based sampler (SDE) $dz_t = v_t(z_t)dt$ $dz_t = f_t(z_t)dt + \sigma_t dW_t$
- Fokker-Planck Equation: marginal density evolution of $p_t(z), t: 0 \rightarrow 1$, for any z (not necessarily for $z = z_t$)

CNF-based sampler (ODE)Diffusion-based sampler (SDE) $\partial_t p_t(z) = -\nabla_z \cdot \left(v_t(z) p_t(z) \right)$ $\partial_t p_t(z) = -\nabla_z \cdot \left(\left[f_t(z) - \frac{\sigma_t^2}{2} \nabla_z \log p_t(z) \right] p_t(z) \right)$

Making both samplers to induce the same density path: $v_t(z) = f_t(z) - \frac{\sigma_t^2}{2} \nabla_z \log p_t(z)$ ("Probability flow ODE")

Converting Between CNF and Diffusion Samplers

• Building the density path by diffusing $z_1 \sim p_1(z_1) \coloneqq \pi(z_1)$:

 $dz_t = \beta_t z_t dt + \sigma \sqrt{2\beta_t} dW_t, t: 1 \to 0$

 $p_t(z_t) = \int N(z_t | \sqrt{1 - \lambda_t} z_1, \sigma^2 \lambda_t I) p_1(z_1) dz_1, \quad \lambda_t = 1 - \exp[-2\int_t^1 \beta_\tau d\tau]$

• Forward ODE/SDE evolution of z_t , $t: 0 \rightarrow 1$, $z_0 \sim p_0(z_0)$

CNF-based sampler (ODE)Diffusion-based sampler (SDE) $dz_t = v_{\phi}(z_t, t)dt$ $dz_t = [\beta_t z_t + 2\sigma^2 \beta_t s_{\phi}(z_t, t)]dt + \sigma \sqrt{2\beta_t} dW_t$

Making both samplers to induce the same density path $\{p_t(z)\}_{t \in [0,1]}$:

 $v_{\phi}(z_t, t) = \beta_t z_t + 2\sigma^2 \beta_t s_{\phi}(z_t, t) - \sigma^2 \beta_t \nabla_z \log p_t(z) \approx \beta_t z_t + \sigma^2 \beta_t s_{\phi}(z_t, t)$

Converting Between CNF and Diffusion Samplers

sample evolution: SDE vs Probability flow ODE, $v_t(z) = f_t(z) - \frac{\sigma_t^2}{2} \nabla_z \log p_t(z)$ $dz_t = f_t(z_t)dt - \sigma_t^2 \nabla_z \log p_t(z_t) + \sigma_t dW_t, t: 1 \to 0$ $dz_t = f_t(z_t)dt + \sigma_t dW_t, t: 0 \to 1$ Z_1 Z_0 $dz_t = v_t(z_t)dt, t: 1 \to 0$ $dz_t = v_t(z_t)dt, t: 0 \rightarrow 1$ SDE Probability Flow ODE

 $\partial_t p_t(z) = -\nabla_z \cdot (v_t(z)p_t(z))$



 $\rightarrow p_1(z) \coloneqq \pi(z)$

Why Can We Construct So Many Different Diffusion/CNF-based Approximations?

• Ultimately, we only care about distribution match $q(z_1) \approx p_1(z_1) \coloneqq \pi(z_1)!$



Unidentifiability: Given $q_{CNF}(z_1) = q_{SDE'}(z_1)$, we cannot show CNF = CNF',

where *CNF*' = probability flow of *SDE*', (unless you make further assumptions)!

Design Principles I

Recall Design Principles I:

Step 1: Specify a divergence $D(q; \pi)$ between variational density q(z) and target $\pi(z)$

Step 2: Find the member of the family $q^* \in Q$ that minimizes the divergence, i.e.,

 $q^* = \operatorname{argmin}_{q \in Q} D(q; \pi)$

Step 3: Use variational density q^* instead of target π in downstream tasks: e.g., $E_{\pi(z)}[F(z)] \approx E_{q^*(z)}[F(z)] \approx \frac{1}{\kappa} \sum_{k=1}^{K} F(z_k)$, $z_k \sim q^*(z)$

Expressivity

- correlations
- Skew and kurtosis
- Multi-modality



Tractability

- Easy and fast to evaluate q(z)
- Fast sampling from its distribution (e.g., to approximate expectations)
- Ease of optimization

Design Principles II

a differentiable and easy-to-evaluate objective $L(q; \pi)$ whose global optimum is $q^* = \pi$ if $\pi \in Q$

Step 1: Specify a divergence $D(q; \pi)$ between variational density q(z) and target $\pi(z)$

Step 2: Find the member of the family $q^* \in Q$ that minimizes the divergence, i.e.,

 $\frac{q^* = argmin_{q \in Q} D(q; \pi)}{q^*} \quad q^* = argmin_{q \in Q} L(q; \pi)$

Step 3: Use variational density q^* instead of target π in downstream tasks: e.g., $E_{\pi(z)}[F(z)] \approx E_{q^*(z)}[F(z)] \approx \frac{1}{K} \sum_{k=1}^{K} F(z_k)$, $z_k \sim q^*(z)$



Tractabilityno need in many casesEasy and fast to evaluate $q(z)^{-1}$ can accept slower sampling speedFast samplingfrom its distribution(e.g., to approximate expectations)

• Ease of optimization and not too slow

Other designs of q distributions



Structured approximations





Implicit approximate posteriors

Particle-based posteriors

Mescheder et al. Adversarial Variational Bayes: Unifying Variational Autoencoders and Generative Adversarial Networks. ICML 2017 Tran et al. Hierarchical Implicit Models and Likelihood-Free Variational Inference. NeurIPS 2017 Li and Turner. Gradient Estimators for Implicit Models. ICLR 2018 Yin and Zhou. Semi-Implicit Variational Inference. ICML 2018



Deriving $\log Z$ estimates \rightarrow an approximate inference method to obtain q



Domke and Sheldon. Divide and Couple: Using Monte Carlo Variational Objectives for Posterior Approximation. NeurIPS 2019 Wainwright and Jordan. Graphical Models, Exponential Families, and Variational Inference. FTML 2008.

Combining approximate inference and sampling



When needing efficient sampling for all p(z|x), $\forall x$: find the optimal proposal distributions via amortized approximate inference

Score-based approximate inference approaches

- Methods for further advances:
 - more expressive families
 - scale to high-dimensions?
- Statistical properties:
 - estimation guarantees?
 - computational vs statistical tradeoffs? (consistency + convergence)
- Applications:
 - Amortized inference?
 - New insights for generative modeling?

Diffusion & Flow based approaches

- Training difficulties: truly simulation-free approach?
 - Not truly simulation-free: requiring (importance) sampling from the density path $\{p_t(z)\}_{t\in[0,1]}$
 - Diffusion: estimating the score
 - CNF with PINN loss: constructing "training data" $q_t(z)$
 - Quality of IS/SMC is crucial to the empirical success!
 - Can we develop other training methods that are truly simulation-free?

Diffusion & Flow based approaches

- "Is it worth it?"
 - Sometimes your scientist friend only cares about **only one** (very difficult) potential function...
 - Maybe yes, in amortized setting: train $q(z|x) \approx p(z|x)$
 - vs (gradient-based) MCMC? Is it really better?
 - E.g. in training energy-based model with contrastive divergence?





Generally not solved yet:

- Better computation & implementation
 - Hardware-aware methods
 - E.g., GPU acceleration with CUDA and/or Triton, etc.



- Robust Implementation of advanced VI in statistics software
 - Flexible *q* distributions
 - Good optimization



ονισια

CUDA

Generally not solved yet:

- Better statistical theory for approximate inference
 - Optimization & statistical properties
 - e.g. (non-)asymptotic behavior, estimation of moments
 - Convergence behavior in optimization?
 - VAE learning for model *p*: bias induced by sub-optimal *Q* family?



Doing probabilistic inference for/with LLMs?

- Speeding-up LLM decoding?
 - E.g., speculative decoding as a "weird way" to do rejection sampling
- LLM "reasoning"?
 - RL-based finetuning: connection with VI/SMC in SSMs
 - "Steered sampling": connections with IS/SMC/MCMC
- Using LLMs to build informative priors?
 - Posterior inference still requires approximations
- LLM "in-context learning" as Bayesian inference?
 - Prediction-centric formulation, still requires approximations





₩Claude

🚫 LLaMA 🗟



Thank You!

Questions? Ask now or contact us via email :)



Diana Cai dcai@flatironinstitute.org



Yingzhen Li yingzhen.li@imperial.ac.uk