## Inference with Scores: Slices, Diffusions and Flows

Yingzhen Li

yingzhen.li@imperial.ac.uk

Feat. Wenbo Gong (University of Cambridge)



#### Bayesian ML / Probability Theory



#### Decision making under uncertainty

Image courtesy of Sebastian Nowozin Re-use of the image for any other purpose is not allowed

# The Central Computation for Inference

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



(For discrete probability measures, integration becomes discrete sum.) (We will discuss continuous variables in the rest of the talk.)

#### **Computation Challenge**

• The central equation for inference:

$$\int F(\theta) \pi(\theta) d\theta$$

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

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



# **Computation Challenge**

• The central equation for inference:

# $\int F(\theta) \pi(\theta) d\theta$

#### "What is the mean of this distribution?"

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



### Approximate Inference

• Central task: approximate  $\pi(\theta)$ 



Approximate distribution design





#### Today's focus

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

min  $Loss(q(\theta), \pi(\theta))$ Image: Contract of the second s

# Approximating the Target Distribution

• Target distribution is usually intractable due to normalising constant:

 $p(\theta|D) = \frac{1}{p(D)} \exp[\log p(\theta) + \log p(D|\theta)]$ 

**Bayesian Posterior** 

$$p(x) = \frac{1}{Z} \exp[-E(x;\theta)]$$

**Energy-based Model** 

• Approximation methods to side-step:

Variational InferenceScore Matching $\min_{q} KL[q||p] = E_q[\log q(x) - \log p(x)]$  $\min_{q} F(q,p) = \frac{1}{2}E_q[||\nabla_x \log q(x) - \nabla_x \log p(x)||_2^2]$ "compare  $\log q(x)$  with  $\log p(x)$ ""compare  $\nabla_x \log q(x)$  with  $\nabla_x \log p(x)$ "Using equiv. ELBO objective,<br/>no need to evaluate  $\log Z$ As  $\nabla_x \log Z = 0$ ,<br/>no need to evaluate  $\log Z$ 

### Comparing KL & Fisher Divergences

• Comparing 2 Gaussians with different mean, same variance ( $x \in \mathbb{R}^d$ ):

 $p(x) = N(x; \mu_p, \sigma^2 I_d)$ 

 $q(x) = N(x; \mu_q, \sigma^2 I_d)$ 

Variational Inference  $KL[q||p] = \frac{1}{2\sigma^2} \|\mu_q - \mu_p\|_2^2$   $F(q, p) = \frac{1}{2\sigma^4} \|\mu_q - \mu_p\|_2^2$ 

Define 
$$\mu_q \coloneqq \mu_p + \sigma \gamma$$
, then  $\|\mu_q - \mu_p\|_2^2 = \|\gamma\|_2^2$ 

#### Comparing KL & Fisher Divergences

• Comparing 2 Gaussians with different mean, same variance ( $x \in \mathbb{R}^d$ ):

high

low

 $p(x) = N(x; \mu_p, \sigma^2 I_d)$ 

Variational Inference





 $q(x) = N(x; \mu_p + \sigma\gamma, \sigma^2 I_d)$ 

Score Matching  $F(q, p) = \frac{1}{2\sigma^2} \|\gamma\|_2^2$ 



Very large if not capturing sharp peaks!

• Fisher divergence: a way to measure score difference:

$$F(p,q) = \frac{1}{2} E_q[\|\nabla_x \log q(x) - \nabla_x \log p(x)\|_2^2]$$
  
Unknown score (only have  $x \sim q(x)$ )



#### Solution: Find a function to directly compute/approximate the score difference!

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

• Stein's identity:

$$E_{p(x)}[\nabla_x \log p(x)^\top f(x) + \nabla_x^\top f(x)] = 0$$

For all  $f(x) \in F$  satisfying  $\lim_{\|x\| \to \infty} p(x)f(x) = 0$ 

• Stein discrepancy:

$$S(q,p) = \sup_{f \in F} E_{q(x)} [\nabla_x \log p(x)^\top f(x) + \nabla_x^\top f(x)]$$

- $q = p \Rightarrow S(q, p)$  (Stein's identity)
- S(q, p) = 0, and if  $f \in F$  is a rich function family  $\Rightarrow q = p$

Gorham and Mackey. Measuring Sample Quality with Stein's Method. NIPS 2015 Liu et al. A Kernelized Stein Discrepancy for Goodness-of-fit Tests. ICML 2016 Chwialkowski et al. A kernel test of goodness of fit. ICML 2016

• Stein discrepancy:

$$S(q,p) = \sup_{f \in F} E_{q(x)} [\nabla_x \log p(x)^\top f(x) + \nabla^\top f(x)],$$

- Estimation requires:
  - Samples  $x \sim q(x) \leftarrow$  enables estimation of  $\nabla_x \log q(x)$  & posterior approximations
  - Score function  $s_p(x)$  evaluated at  $x \leftarrow \text{enables learning of } p(x)$  as an EBM

$$p(x) = \frac{1}{Z} \exp[-E_{\theta}(x)]$$
$$\nabla_{x} \log p(x) = -\nabla_{x} E_{\theta}(x) - \overline{\nabla_{x} \log Z}$$
$$= 0$$

• Stein discrepancy:

$$S(q,p) = \sup_{f \in F} E_{q(x)} [\nabla_x \log p(x)^\top f(x) + \nabla^\top f(x)],$$

- Selecting test functions  $f \in F$ :
  - As  $L_2$  integrable functions:
    - $f^*(x) \propto \nabla_x \log p(x) \nabla_x \log q(x)$  (intractable)
    - Stein discrepancy  $\Leftrightarrow$  Fisher divergence (integration by parts)
    - Optimal test function approximated by a neural network  $f_{\phi}(x): \mathbb{R}^D \to \mathbb{R}^D$

Ranganath et al. Operator Variational Inference. NIPS 2016

Grathwohl et al. Learning the Stein Discrepancy for Training and Evaluating Energy-Based Models without Sampling. ICML 2020

• Stein discrepancy:

$$S(q,p) = \sup_{f \in F} E_{q(x)} [\nabla_x \log p(x)^\top f(x) + \nabla^\top f(x)],$$

- Selecting test functions  $f \in F$ :
  - As functions in a unit ball of an RKHS:
    - Closed-form optimal test function as "smoothed" score difference
    - Optimal test function can be computed (integration by parts)
    - Test power depends on the choice of kernel

• Stein discrepancy:

$$S(q,p) = \sup_{f \in F} E_{q(x)} [\nabla_x \log p(x)^\top f(x) + \nabla^\top f(x)],$$

- Curse-of-dimensionality problem:
  - $f \in F$  as  $L_2$  integrable functions:
    - Optimizing  $f_{\phi}(x): \mathbb{R}^D \to \mathbb{R}^D$  is challenging in high dimensions
  - $f \in F$  as functions in a unit ball of an RKHS:
    - Open question of kernel choice in high dimensions

# Sliced Stein Discrepancy

• 1<sup>st</sup> idea: find the projection direction where scores differ the most (on average)!



coordinate system with basis

 $I_d = \begin{bmatrix} 1 \ 0 \ \cdots \ 0 \\ \dots \\ 0 \ 0 \ \dots \ 1 \end{bmatrix}$ 

 $\nabla_x \log p(x) - \nabla_x \log q(x)$ 



coordinate system with basis

$$O_r^{\top} = \begin{bmatrix} r_1^{\top} \\ \dots \\ r_d^{\top} \end{bmatrix}$$

 $O_r^{\top}(\nabla_x \log p(x) - \nabla_x \log q(x))$ 

One best projection is enough!

 $\gamma_*$ 

 $r_*^{\top}(\nabla_x \log p(x) - \nabla_x \log q(x))$ 

Gong et al. Sliced Kernelized Stein Discrepancy. ICLR 2021 Gong et al. Active Slices for Sliced Stein Discrepancy. ICML 2021

# Sliced Stein Discrepancy

• 1<sup>st</sup> idea: find the projection direction where scores differ the most (on average)!



Gong et al. Sliced Kernelized Stein Discrepancy. ICLR 2021 Gong et al. Active Slices for Sliced Stein Discrepancy. ICML 2021

# Sliced Stein Discrepancy

• 2<sup>nd</sup> idea: apply Radon transform

 $f_r: \mathbb{R}^d \to \mathbb{R} \quad \Leftrightarrow \quad f_{rg}: \mathbb{R} \to \mathbb{R}, \, \forall g \in S^{d-1}$ 

• ... again, pick best g

Resulting Stein discrepancy: maxSSD-rg

 $\sup_{r,g,f_{rg}\in F_{rg}} E_{q(x)}[r^{\top}\nabla_{x}\log p(x)f_{rg}(x^{T}g) + r^{T}g\nabla_{x^{T}g}f_{rg}(x^{T}g)]$ 

 $f_{rg}: R \to R$ 



Gong et al. Sliced Kernelized Stein Discrepancy. ICLR 2021 Gong et al. Active Slices for Sliced Stein Discrepancy. ICML 2021

# Goodness-of-fit Test with SSD

#### Goodness-of-fit test:

- Both *p* and *q* are Gaussian-RBMs
- q weights as perturbed from p weights

- Better test power:
  - at perturbation level 0.01, rejection rate 95% (Active slice) vs 45%(GO)
- Significant speed-up:
  - 0.04s (Active slice) vs 10.15s (GO)
  - 254x times faster



#### Selecting SGHMC Step-sizes with SD



#### From Global to Local Coordinates



Stein discrepancy (SD) Global Coordinate system  $I_d$ 

Sliced Stein discrepancy (SSD) Global Coordinate system  $O_r$  Diffusion score matching & SD Local Coordinate system

Barp et al. Minimum Stein Discrepancy Estimators. NeurIPS 2019 Gong and Li. Interpreting Diffusion Score Matching using Normalizing Flows. ICML 2021 INNF+ workshop

#### Fitting Student-t with Score Matching



$$\operatorname{argmin}_{\theta} \mathcal{F}(\boldsymbol{q}_{\text{data}}, \boldsymbol{p}_{\theta}) = \operatorname{argmin}_{\theta} \mathcal{L}(\boldsymbol{\theta})$$

### Fitting Student-t with Score Matching

 $q_{\text{data}}(\boldsymbol{x}) = \text{Student-t}(0, 0.3)$  $p_{\theta}(\boldsymbol{x}) = \text{Student-t}(\theta, 0.3)$ SM and DSM Loss for Student-t distribution at different  $\theta$ Score Matching Valid region Loss 2 0 -2 -4-2 -4 2 -60 6 4 Mean  $\theta$ 

- $\theta$  initialized in valid region  $\Rightarrow$  optimisation converges to  $\theta^* = 0$  (global optimum)
- Otherwise, optimisation converges to  $\pm \infty$  (local optimum)

#### Fitting Student-t with Score Matching





How to interpret & choose m(x)?

### Interpreting DSM using Flows



#### Similar idea applies to Diffusion Stein Discrepancy

#### Better Flow Design

Flow (Barp et al.)

$$oldsymbol{m}oldsymbol{x}ig) \quad oldsymbol{m}oldsymbol{x}) = 1 + rac{(oldsymbol{x} - heta)^2}{0.6}$$

#### **Gaussian Flow**

 $\boldsymbol{m}(\boldsymbol{x}) = (\nabla_{\boldsymbol{x}} T(\boldsymbol{x}))^{-1}$ 

where  $T(\cdot)$  transforms  $p_{\theta}$  to standard Gaussian  $T(\cdot)$  is composition of  $\text{CDF}_{p_{\theta}}$  and  $\text{CDF}_{\mathcal{N}}^{-1}$ 



Faster convergence with Gaussian Flow compared to Flow by Barp et al. 2019

#### The Name "Diffusion"

- Stein operator ⇒ Stein discrepancy:
  - Find Stein operator  $A_p[f]$  such that  $E_{p(x)}[A_p[f](x)] = 0$  for  $\forall f \in F$
  - Construct the corresponding Stein discrepancy  $S(q,p) = \sup_{f \in F} E_{q(x)}[A_p[f](x)]$

$$dx = b(x)dt + \sqrt{2D(x)}dW(t)$$
$$b(x) = [D(x) + Q(x)]\nabla_x \log p(x) + \Gamma(x)$$

generator method

$$A_p[f](x) = \frac{1}{p(x)} \langle \nabla, p(x) (D(x) + Q(x)) f(x) \rangle$$

Complete SG-MCMC recipe Ito diffusion SDE Diffusion Stein discrepancy with m(x) = D(x) + Q(x)

Gorham et al. Measuring Sample Quality with Diffusions. Annuals of Applied Probability 2019 Ma et al. A Complete Recipe for Stochastic Gradient MCMC. NIPS 2015

#### Conclusion

- Score-matching & Stein discrepancy as promising alternatives
  - Still facing challenges in high dimensions
- Ideas from flows can help improve Stein's method
- How about the other direction?

Thank you!