Real hero behind many relevant works:

Wenbo Gong (MSR)



Make Stein's method "great again" for generative modelling?

Yingzhen Li

yingzhen.li@imperial.ac.uk



Training generative models with "gradient approximation"

 $\min_{\phi} L(\phi)$ with intractable $L(\phi)$ Loss approximation: Gradient approximation: Find $\tilde{L} \approx L$ then $\min_{\phi} \tilde{L}(\phi)$ Find $\tilde{g}(\phi) \approx \nabla_{\phi} L$ then update $\phi \leftarrow \phi - \gamma \tilde{g}(\phi)$ true loss true minimum true gradient true minimum approx. loss ★ approx. loss minima approx. gradient

Example: applications in variational inference

Variational lower-bound: assume $\pmb{z} \sim \pmb{q_\phi} \Leftrightarrow \pmb{\epsilon} \sim \pi(\pmb{\epsilon}), \pmb{z} = \pmb{f_\phi}(\pmb{\epsilon}, \pmb{x})$

$$egin{aligned} \mathcal{L}_{\mathsf{VI}}(q_{\phi}) &= \mathbb{E}_{q} \left[\log p(oldsymbol{x},oldsymbol{z})
ight] + \mathbb{H}[q_{\phi}(oldsymbol{z}|oldsymbol{x})] \ &= \mathbb{E}_{\pi} \left[\log p(oldsymbol{x},oldsymbol{f}_{\phi}(\epsilon,oldsymbol{x}))] + \mathbb{H}[q_{\phi}(oldsymbol{z}|oldsymbol{x})] \ // \ reparam. \ trick \end{aligned}$$

If you use gradient descent for optimisation, then you only need gradients! The gradient of the variational lower-bound:

$$\nabla_{\phi} \mathcal{L}_{\mathsf{VI}}(q_{\phi}) = \mathbb{E}_{\pi} \left[\nabla_{f} \log p(\boldsymbol{x}, f_{\phi}(\boldsymbol{\epsilon}, \boldsymbol{x}))^{\mathsf{T}} \nabla_{\phi} f_{\phi}(\boldsymbol{\epsilon}, \boldsymbol{x}) \right] + \nabla_{\phi} \mathbb{H}[q_{\phi}(\boldsymbol{z} | \boldsymbol{x})]$$

The gradient of the entropy term:

$$\nabla_{\phi} \mathbb{H}[q_{\phi}(\boldsymbol{z}|\boldsymbol{x})] = -\mathbb{E}_{\pi}\left[\nabla_{f} \log q(f_{\phi}(\boldsymbol{\epsilon}, \boldsymbol{x})|\boldsymbol{x})^{\mathsf{T}} \nabla_{\phi} f_{\phi}(\boldsymbol{\epsilon}, \boldsymbol{x})\right] - \underbrace{\mathbb{E}_{q} \left[\nabla_{\phi} \log q_{\phi}(\boldsymbol{z}|\boldsymbol{x})\right]}_{\Phi}$$

this term is 0

It remains to approximate $\nabla_z \log q(z|x)!$

Li and Turner. Gradient Estimators for Implicit Models. ICLR 2018

Goal: approximate $\nabla_x \log q(x)$ for a given distribution q(x)Main idea: invert Stein's identity:

 $\mathbb{E}_q[\boldsymbol{h}(\boldsymbol{x})\nabla_{\boldsymbol{x}}\log q(\boldsymbol{x})^{\mathsf{T}}+\nabla_{\boldsymbol{x}}\boldsymbol{h}(\boldsymbol{x})]=\boldsymbol{0}$

1. Monte Carlo (MC) approximation to Stein's identity:

$$\frac{1}{K}\sum_{k=1}^{K} -\boldsymbol{h}(\boldsymbol{x}^{k}) \nabla_{\boldsymbol{x}^{k}} \log q(\boldsymbol{x}^{k})^{\mathsf{T}} + \mathsf{err} = \frac{1}{K}\sum_{k=1}^{K} \nabla_{\boldsymbol{x}^{k}} \boldsymbol{h}(\boldsymbol{x}^{k}), \quad \boldsymbol{x}^{k} \sim q(\boldsymbol{x}^{k}),$$

2. Rewrite the MC equations in matrix forms: denoting

$$\mathbf{H} = \left(\mathbf{h}(\mathbf{x}^{1}), \cdots, \mathbf{h}(\mathbf{x}^{K}) \right), \quad \overline{\nabla_{\mathbf{x}}\mathbf{h}} = \frac{1}{K} \sum_{k=1}^{K} \nabla_{\mathbf{x}^{k}} \mathbf{h}(\mathbf{x}^{k}),$$
$$\mathbf{G} := \left(\nabla_{\mathbf{x}^{1}} \log q(\mathbf{x}^{1}), \cdots, \nabla_{\mathbf{x}^{K}} \log q(\mathbf{x}^{K}) \right)^{\mathsf{T}},$$
Then $-\frac{1}{K} \mathbf{H} \mathbf{G} + \operatorname{err} = \overline{\nabla_{\mathbf{x}}\mathbf{h}}.$

Goal: approximate $\nabla_x \log q(x)$ for a given distribution q(x)Main idea: invert Stein's identity:

 $\mathbb{E}_q[\boldsymbol{h}(\boldsymbol{x})\nabla_{\boldsymbol{x}}\log q(\boldsymbol{x})^{\mathsf{T}}+\nabla_{\boldsymbol{x}}\boldsymbol{h}(\boldsymbol{x})]=\boldsymbol{0}$

Matrix form (MC): $-\frac{1}{K}$ **HG** + err = $\overline{\nabla_x h}$.

3. Now solve a ridge regression problem:

$$\hat{\mathbf{G}}_{V}^{\mathsf{Stein}} := \argmin_{\hat{\mathbf{G}} \in \mathbb{R}^{K \times d}} ||\overline{\nabla_{\boldsymbol{x}} \boldsymbol{h}} + \frac{1}{K} \mathbf{H} \hat{\mathbf{G}}||_{F}^{2} + \frac{\eta}{K^{2}} ||\hat{\mathbf{G}}||_{F}^{2},$$

Analytic solution: $\hat{\mathbf{G}}_{V}^{\text{Stein}} = -(\mathbf{K} + \eta \mathbf{I})^{-1} \langle \nabla, \mathbf{K} \rangle$, with $\mathbf{K} := \mathbf{H}^{\mathsf{T}} \mathbf{H}, \quad \mathbf{K}_{ij} = \mathcal{K}(\mathbf{x}^{i}, \mathbf{x}^{j}) := \mathbf{h}(\mathbf{x}^{i})^{\mathsf{T}} \mathbf{h}(\mathbf{x}^{j}),$ $\langle \nabla, \mathbf{K} \rangle := \mathbf{K} \mathbf{H}^{\mathsf{T}} \overline{\nabla_{\mathbf{x}} \mathbf{h}}, \quad \langle \nabla, \mathbf{K} \rangle_{ij} = \sum_{k=1}^{K} \nabla_{\mathbf{x}_{i}^{k}} \mathcal{K}(\mathbf{x}^{i}, \mathbf{x}^{k}).$

Kernelized Stein Discrepancy:

$$\begin{split} \mathcal{S}^2(q,\hat{q}) &= \mathbb{E}_{\mathbf{x},\mathbf{x}'\sim q} \left[\hat{\mathbf{g}}(\mathbf{x})^\mathsf{T} \mathcal{K}_{\mathbf{x}\mathbf{x}'} \hat{\mathbf{g}}(\mathbf{x}') + \hat{\mathbf{g}}(\mathbf{x})^\mathsf{T} \nabla_{\mathbf{x}'} \mathcal{K}_{\mathbf{x}\mathbf{x}'} + \nabla_{\mathbf{x}} \mathcal{K}_{\mathbf{x}\mathbf{x}'}^\mathsf{T} \hat{\mathbf{g}}(\mathbf{x}') + \mathsf{Tr}(\nabla_{\mathbf{x},\mathbf{x}'} \mathcal{K}_{\mathbf{x}\mathbf{x}'}) \right], \\ \mathbf{g}(\mathbf{x}) &= \nabla_{\mathbf{x}} \log q(\mathbf{x}), \quad \hat{\mathbf{g}}(\mathbf{x}) = \nabla_{\mathbf{x}} \log \hat{q}(\mathbf{x}), \quad \mathcal{K}_{\mathbf{x}\mathbf{x}'} = \mathcal{K}(\mathbf{x},\mathbf{x}'). \end{split}$$

One can show that the V-statistic of KSD is

$$\mathcal{S}_{V}^{2}(q,\hat{q}) = \frac{1}{K^{2}} \operatorname{Tr}(\hat{\mathbf{G}}^{\mathsf{T}}\mathbf{K}\hat{\mathbf{G}} + 2\hat{\mathbf{G}}^{\mathsf{T}}\langle \nabla, \mathbf{K} \rangle) + C$$

This means

$$\hat{\mathbf{G}}_{V}^{Stein} = \operatorname*{arg\,min}_{\hat{\mathbf{G}} \in \mathbb{R}^{K \times d}} \quad \mathcal{S}_{V}^{2}(q, \hat{q}) + \frac{\eta}{K^{2}} ||\hat{\mathbf{G}}||_{F}^{2}$$

Example: Training an implicit generative model p_{θ} by minimising $KL[p_{\theta} || p_{data}]$:



(defined like GAN's generator)

The rise of score-based generative models



Would Stein's method ever work for generative models?

What do I mean by "work": which one is the generator trained using Stein's method?



Spot the connections & differences



Solving common issues

- Score function is an $R^d \rightarrow R^d$ mapping
 - Issues in estimating it when d is large, e.g., $~d \geq 100~$

Solution: slicing



• 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:
 - $f^*(x) \propto \nabla_x \log p(x) \nabla_x \log q(x)$ (intractable)
 - 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

• 1st 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!

 γ_*

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

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

• 1st 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

• 2nd 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

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

 $f_{rg}: R \to R$





ICA Model Learning

ICA generative process:

 $z \sim Lap(0,1), x = Wz, \log p(x) = \log p_z(W^{-1}x) + C$ With normalizing constant *C*.

Training/evaluation setup:

- Generate dataset \boldsymbol{x} by ICA with ground truth \boldsymbol{W}_t
- Model $p(\mathbf{x})$ is an ICA with random initialized \mathbf{W} .
- Train p(x) to match dataset x.
- Evaluate with log likelihood on test data.

ICA Model Learning

Method	Dimension						
	D = 10	D = 20	D = 40	D = 60	D = 80	D = 100	D = 200
KSD	-10.23	-15.98	-34.50	-56.87	-86.09	-116.51	-329.49
LSD	-10.42	-14.54	-17.16	-15.05	-12.39	-5.49	46.63
maxSKSD-9	-10.45	-14.50	-17.28	-15.70	-11.91	-4.21	47.72



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

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

Solving common issues

Disconnected modes

Solution: smoothing approximations



Data scores	Estimated scores				
/ / / / / / / / / / / / / / / / / / /	MARAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA				
1 ************************************	111111111111111111111111111111111111111				
11	/////				
11/2	///////////////////////////////////////				
111/	///////////////////////////////////////				
1111	///////////////////////////////////////				
The second secon	CONTRACTOR CONTRACTOR STATES				
* * * * * * * * * * * * * * * * * * * *	*******				
********	***********************				
· · · · · · · · · · · · · · · · · · ·	******************				
	111111111111111111111111111111111111111				
	111111111111				
[· · · · · ·] · · · · · · · · · · · ·					

• • • • • • • • • • • • • • • • • • •	* · · · · · · · · · · · · · · · · · · ·				
1,1 / * * * * * * * * * * * * * * * * * *	111111111111111111111111111111111111111				
Kin in right a second of 1	KILLING CONTRACTOR CONTRACTOR				

Example:

$$p(x) = \pi_1 p_1(x) + (1 - \pi_1) p_2(x)$$

$$\nabla \log p(x) = \frac{\pi_1 p_1(x)}{\pi_1 p_1(x) + (1 - \pi_1) p_2(x)} \nabla \log p_1(x)$$

$$+ \frac{\pi_2 p_2(x)}{\pi_1 p_1(x) + (1 - \pi_1) p_2(x)} \nabla \log p_2(x)$$

Adding noise & Tempering: "temperature-matched" KSD:



Wenliang and Kanagawa. Blindness of score-based methods to isolated components and mixing proportions. arXiv:2008.10087

Denoising Autoencoder objective:

$$L(r_{\sigma}) = E_{p(x)p_{\sigma}(x'|x)}[\|r_{\sigma}(x') - x\|_{2}^{2}]$$

Performing second-order Taylor expansion:

$$r_{\sigma}(x + \sigma\epsilon) = r_{\sigma}(x) + \sigma \nabla_{x} r_{\sigma}(x)\epsilon + O(\sigma^{2}\epsilon^{2})$$

Rewriting the DAE objective:

$$L(r_{\sigma}) = E_{p(x)} [\|r_{\sigma}(x) - x\|_{2}^{2} + \sigma^{2} \|\nabla_{x} r_{\sigma}(x)\|_{F}^{2}] + O(\sigma^{2})$$

Alain and Bengio (2014). What regularized auto-encoders learn from the data-generating distribution. JMLR. Unpublished result with Wenbo Gong.



Denoising Autoencoder objective when σ is small:

$$L(r_{\sigma}) = E_{p(x)} [\| r_{\sigma}(x) - x \|_{2}^{2} + \sigma^{2} \| \nabla_{x} r_{\sigma}(x) \|_{F}^{2}]$$

Consider finding the functional gradient in an RKHS H^d with kernel k:

$$F(f) \coloneqq L(r_{\sigma} + f), \quad \nabla_{f}F(f) \coloneqq \lim_{\epsilon \to 0} \frac{F(f + \epsilon g) - F(f)}{\epsilon}$$

We can show:

$$\nabla_f F(f)|_{f=0} = 2\sigma^2 E_{p(x)}[(r_\sigma(x) - x)k(x, \cdot) + \nabla_x k(x, \cdot)]$$

Define
$$q(x)$$
 such that $\nabla_x \log q(x) = \frac{r_{\sigma}(x) - x}{\sigma^2} \implies \left\| \nabla_f F(f) \right\|_{f=0} \right\|_{H^d}^2 = 4\sigma^4 KSD(p,q)$

Corrupted input

Denoising Autoencoder objective when σ is small:

$$L(r_{\sigma}) = E_{p(x)} [\| r_{\sigma}(x) - x \|_{2}^{2} + \sigma^{2} \| \nabla_{x} r_{\sigma}(x) \|_{F}^{2}]$$

Functional gradient in an RKHS H^d with kernel k:

$$\nabla_f F(f)|_{f=0} = 2\sigma^2 E_{p(x)}[(r_\sigma(x) - x)k(x, \cdot) + \nabla_x k(x, \cdot)]$$

With
$$\nabla_x \log q(x) = \frac{r_\sigma(x) - x}{\sigma^2}$$
, $\|\nabla_f F(f)\|_{f=0}\|_{H^d}^2 = 4\sigma^4 KSD(p,q)$

- Can try to find DAE updates in an RKHS following the functional gradient
- KSD captures the RKHS norm of the functional gradient
- Finding optimal DAE $r_{\sigma}(\cdot) \Leftrightarrow KSD(p,q) = 0$

Corrupted input





Big problem: finding a suitable test function family!

Solving specific issues?

- Kernel methods ever suitable?
- Answer (?): Deep Kernel learning

 $\min_{\theta} \max_{\phi} KSD_{\phi}(p_{\theta}, p_{data})?$ $K_{\phi}(x, y) = K_{RBF}(f_{\phi}(x), f_{\phi}(y))$



MMD-GAN generated CelebA images

Wilson et al. Deep Kernel Learning. AISTATS 2016

Sutherland et al. Generative Models and Model Criticism via Optimized Maximum Mean Discrepancy. ICLR 2017 Li et al. MMD GAN: Towards Deeper Understanding of Moment Matching Network. NeurIPS 2017

Solving specific issues?

- Adversarial training for multi-dimensional test functions?
- Developments for GANs: lots of stabilisation tricks
 - Label smoothing
 - Regularising the norm of the discriminator function
 - Neural network normalisation layers
 - Designs of learning rate schedules for the generator & discriminator
 - .

Are these tricks directly transferable to Learned Stein discrepancies?

Salimans et al. Improved Techniques for Training GANs. NIPS 2016 Gulrajani et al. Improved Training of Wasserstein GANs. NIPS 2017 Miyato et al. Spectral Normalization for Generative Adversarial Networks. ICLR 2018 Heusel et al. GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium. NIPS 2017

An open question

- Does estimation accuracy of the score matters a lot always?
 - Important for later stages

(Text-conditioned) Generative Denoising Process





data size & model size

Thanks to my awesome collaborators:





...and the future MRes student.



Wenbo Gong

Richard E Turner Jose Miguel Hernandez-Lobato

(for believing in me in this topic)