

Gradient Estimators for Implicit Models

Yingzhen Li

University of Cambridge

Joint work with Rich Turner (Cambridge)

Special thanks to Qiang Liu (Dartmouth → UT Austin)

Based on Chapter 5,6 of my (draft) thesis

arXiv 1705.07107

Tractability in approximate inference

- Bayesian inference: integrating out the unobserved variables in your model
 - latent variables in a generative model
 - weight matrices in a Bayesian neural network
- “We do approximate inference because exact inference is **intractable**.”

Tractability in approximate inference

- Bayesian inference: integrating out the unobserved variables in your model
 - latent variables in a generative model
 - weight matrices in a Bayesian neural network
- “We do approximate inference because exact inference is **intractable**.”

What does “**tractability**” really mean for an *approximate* inference algorithm?

- Bayes Rule

$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{z})p(\mathbf{x}|\mathbf{z})}{p(\mathbf{x})}$$

- **Inference:** given some function $F(\mathbf{z})$ in interest, want $\mathbb{E}_{p(\mathbf{z}|\mathbf{x})} [F(\mathbf{z})]$
 - predictive distribution $p(\mathbf{y}|\mathbf{x}, \mathcal{D}) = \mathbb{E}_{p(\mathbf{z}|\mathcal{D})} [p(\mathbf{y}|\mathbf{x}, \mathbf{z})]$
 - evaluate posterior $p(\mathbf{z} \in A|\mathbf{x}) = \mathbb{E}_{p(\mathbf{z}|\mathbf{x})} [\delta_A]$
- In this talk we assume $F(\mathbf{z})$ is cheap to compute given \mathbf{z}

Tractability in approximate inference

- In most of the time we cannot compute $p(\mathbf{z}|\mathbf{x})$ efficiently
- **Approximate inference**: find $q(\mathbf{z}|\mathbf{x})$ in some family \mathcal{Q} such that $q(\mathbf{z}|\mathbf{x}) \approx p(\mathbf{z}|\mathbf{x})$
- In inference time: Monte Carlo approximation:

$$\mathbb{E}_{p(\mathbf{z}|\mathbf{x})} [F(\mathbf{z})] \approx \frac{1}{K} \sum_{k=1}^K F(\mathbf{z}^k), \quad \mathbf{z}^k \sim q(\mathbf{z}|\mathbf{x})$$

Tractability in approximate inference

- In most of the time we cannot compute $p(\mathbf{z}|\mathbf{x})$ efficiently
- **Approximate inference**: find $q(\mathbf{z}|\mathbf{x})$ in some family \mathcal{Q} such that $q(\mathbf{z}|\mathbf{x}) \approx p(\mathbf{z}|\mathbf{x})$
- In inference time: Monte Carlo approximation:

$$\mathbb{E}_{p(\mathbf{z}|\mathbf{x})} [F(\mathbf{z})] \approx \frac{1}{K} \sum_{k=1}^K F(\mathbf{z}^k), \quad \mathbf{z}^k \sim q(\mathbf{z}|\mathbf{x})$$

Tractability requirement: **fast sampling** from q

Tractability in approximate inference

Optimisation-based methods, e.g. **variational inference**:

- Optimise a (usually parametric) q distribution to approximate the exact posterior

$$q^*(\mathbf{z}|\mathbf{x}) = \arg \min_{q \in \mathcal{Q}} \text{KL}[q(\mathbf{z}|\mathbf{x})||p(\mathbf{z}|\mathbf{x})] = \arg \max_{q \in \mathcal{Q}} \mathbb{E}_q[\log p(\mathbf{z}, \mathbf{x})] + \mathbb{H}[q(\mathbf{z}|\mathbf{x})]$$

- When q or p is complicated, usually approximate the expectation by Monte Carlo

$$\mathcal{L}_{\text{VI}}^{\text{MC}}(q) = \frac{1}{K} \sum_{k=1}^K \log p(\mathbf{x}, \mathbf{z}^k) - \log q(\mathbf{z}^k|\mathbf{x}), \quad \mathbf{z}^k \sim q(\mathbf{z}|\mathbf{x})$$

- With Monte Carlo approximation methods, inference is done by

$$\mathbb{E}_{p(\mathbf{z}|\mathbf{x})} [F(\mathbf{z})] \approx \frac{1}{K} \sum_{k=1}^K F(\mathbf{z}^k), \quad \mathbf{z}^k \sim q(\mathbf{z}|\mathbf{x})$$

Tractability in approximate inference

Optimisation-based methods, e.g. **variational inference**:

- Optimise a (usually parametric) q distribution to approximate the exact posterior

$$q^*(\mathbf{z}|\mathbf{x}) = \arg \min_{q \in \mathcal{Q}} \text{KL}[q(\mathbf{z}|\mathbf{x})||p(\mathbf{z}|\mathbf{x})] = \arg \max_{q \in \mathcal{Q}} \mathbb{E}_q[\log p(\mathbf{z}, \mathbf{x})] + \mathbb{H}[q(\mathbf{z}|\mathbf{x})]$$

- When q or p is complicated, usually approximate the expectation by Monte Carlo

$$\mathcal{L}_{\text{VI}}^{\text{MC}}(q) = \frac{1}{K} \sum_{k=1}^K \log p(\mathbf{x}, \mathbf{z}^k) - \log q(\mathbf{z}^k|\mathbf{x}), \quad \mathbf{z}^k \sim q(\mathbf{z}|\mathbf{x})$$

- With Monte Carlo approximation methods, inference is done by

$$\mathbb{E}_{p(\mathbf{z}|\mathbf{x})} [F(\mathbf{z})] \approx \frac{1}{K} \sum_{k=1}^K F(\mathbf{z}^k), \quad \mathbf{z}^k \sim q(\mathbf{z}|\mathbf{x})$$

Tractability requirement: **fast sampling**
and **fast density (gradient) evaluation** (only for optimisation)

Is it necessary to evaluate the (approximate) posterior density?

Three reasons why I think it is not necessary:

- if yes, might restrict the approximation accuracy
- if yes, visualising distributions in high dimensions is still an open research question
- most importantly, MC integration does not require density evaluation

Can we design efficient approximate inference algorithms that enables **fast** inference, **without adding more requirements to q ?**

Why this research problem is interesting:

- Having the best from both MCMC and VI
- Allowing exciting new applications

VI

- Need fast density (ratio) evaluation
- Less accurate
- Faster inference
- Easy to amortise (memory efficient)

MCMC

- Just need to do sampling
- Very accurate
- Need large T thus slower
- Not re-usable when p is updated

We want to have the best from both worlds!

Wild approximate inference: Why

Meta learning for approximate inference:

- Currently we handcraft MCMC algorithms and/or approximate inference optimisation objectives
- Can we learn them?

Wild approximate inference: How

We have seen/described/developed 4 categories of approaches:

- Variational lower-bound approximation (based on density ratio estimation)
Li and Liu (2016), Karaletsos (2016); Huszár (2017); Tran et al. (2017); Mescheder et al. (2017); Shi et al. (2017)
- Alternative objectives other than minimising KL
Ranganath et al. (2016); Liu and Feng (2016)
- Amortising deterministic/stochastic dynamics
Wang and Liu (2016); Li, Turner and Liu (2017); Chen et al. (2017); Pu et al. (2017)
- Gradient approximations (this talk)
Huszár (2017); Li and Turner (2017)

Also see Titsias (2017)

Alternative idea: approximate the gradient

Variational lower-bound: assume $\mathbf{z} \sim q_\phi \Leftrightarrow \epsilon \sim \pi(\epsilon), \mathbf{z} = \mathbf{f}_\phi(\epsilon)$

$$\mathcal{L}_{\text{VI}}(q_\phi) = \mathbb{E}_\pi [\log p(\mathbf{x}, \mathbf{f}_\phi(\epsilon, \mathbf{x}))] + \mathbb{H}[q(\mathbf{z}|\mathbf{x})]$$

During optimisation we only care about the gradients!

The gradient of the variational lower-bound:

$$\nabla_\phi \mathcal{L}_{\text{VI}}(q_\phi) = \mathbb{E}_\pi [\nabla_{\mathbf{f}} \log p(\mathbf{x}, \mathbf{f}_\phi(\epsilon, \mathbf{x}))^\top \nabla_\phi \mathbf{f}_\phi(\epsilon, \mathbf{x})] + \nabla_\phi \mathbb{H}[q(\mathbf{z}|\mathbf{x})]$$

Alternative idea: approximate the gradient

Variational lower-bound: assume $\mathbf{z} \sim q_\phi \Leftrightarrow \epsilon \sim \pi(\epsilon), \mathbf{z} = \mathbf{f}_\phi(\epsilon)$

$$\mathcal{L}_{\text{VI}}(q_\phi) = \mathbb{E}_\pi [\log p(\mathbf{x}, \mathbf{f}_\phi(\epsilon, \mathbf{x}))] + \mathbb{H}[q(\mathbf{z}|\mathbf{x})]$$

During optimisation we only care about the gradients!

The gradient of the variational lower-bound:

$$\nabla_\phi \mathcal{L}_{\text{VI}}(q_\phi) = \mathbb{E}_\pi [\nabla_{\mathbf{f}} \log p(\mathbf{x}, \mathbf{f}_\phi(\epsilon, \mathbf{x}))^\top \nabla_\phi \mathbf{f}_\phi(\epsilon, \mathbf{x})] + \nabla_\phi \mathbb{H}[q(\mathbf{z}|\mathbf{x})]$$

The gradient of the entropy term:

$$\nabla_\phi \mathbb{H}[q(\mathbf{z}|\mathbf{x})] = -\mathbb{E}_\pi [\nabla_{\mathbf{f}} \log q(\mathbf{f}_\phi(\epsilon, \mathbf{x})|\mathbf{x})^\top \nabla_\phi \mathbf{f}_\phi(\epsilon, \mathbf{x})] - \cancel{\mathbb{E}_q [\nabla_\phi \log q_\phi(\mathbf{z}|\mathbf{x})]}$$

this term is 0

It remains to approximate $\nabla_{\mathbf{z}} \log q(\mathbf{z}|\mathbf{x})!$
(in a cheap way, don't want double-loop)

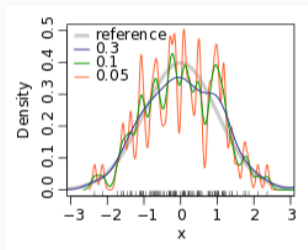
Gradient estimators (kernel based)

KDE plug-in gradient estimator for $\nabla_{\mathbf{x}} \log q(\mathbf{x})$ for $\mathbf{x} \in \mathbb{R}^d$:

- first approximate $q(\mathbf{x})$ using kernel density estimator $\hat{q}(\mathbf{x})$:

$$\hat{q}(\mathbf{x}) = \frac{1}{K} \sum_{k=1}^K \mathcal{K}(\mathbf{x}, \mathbf{x}^k), \quad \mathbf{x}^k \sim q(\mathbf{x})$$

- then approximate $\nabla_{\mathbf{x}} \log q(\mathbf{x}) \approx \nabla_{\mathbf{x}} \log \hat{q}(\mathbf{x})$



Gradient estimators (kernel based)

Score matching gradient estimators: find $\hat{\mathbf{g}}(\mathbf{x})$ to minimise the ℓ_2 error

$$\mathcal{F}(\hat{\mathbf{g}}) := \mathbb{E}_q [\|\hat{\mathbf{g}}(\mathbf{x}) - \nabla_{\mathbf{x}} \log q(\mathbf{x})\|_2^2]$$

Using **integration by parts** we can rewrite: (Hyvärinen 2005)

$$\mathcal{F}(\hat{\mathbf{g}}) = \mathbb{E}_q \left[\|\hat{\mathbf{g}}(\mathbf{x})\|_2^2 + 2 \sum_{j=1}^d \nabla_{x_j} \hat{g}_j(\mathbf{x}) \right] + C$$

Sasaki et al. (2014) and Strathmann et al. (2015): define

$$\hat{\mathbf{g}}(\mathbf{x}) = \sum_{k=1}^K a_k \nabla_{\mathbf{x}} \mathcal{K}(\mathbf{x}, \mathbf{x}^k), \quad \mathbf{x}^k \sim q(\mathbf{x})$$

and find the best $\mathbf{a} = (a_1, \dots, a_K)$ by minimising the ℓ_2 error.

Stein gradient estimator (kernel based)

Define $\mathbf{h}(\mathbf{x})$: a (column vector) test function satisfying the **boundary condition**

$$\lim_{\mathbf{x} \rightarrow \infty} q(\mathbf{x})\mathbf{h}(\mathbf{x}) = \mathbf{0}.$$

Then we can derive **Stein's identity** using **integration by parts**:

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

Invert Stein's identity to obtain $\nabla_{\mathbf{x}} \log q(\mathbf{x})!$

Stein gradient estimator (kernel based)

Main idea: invert Stein's identity:

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

1. MC approximation to Stein's identity:

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

Stein gradient estimator (kernel based)

Main idea: **invert Stein's identity**:

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

1. MC approximation to Stein's identity:

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

2. Rewrite the MC equations in matrix forms: denoting

$$\mathbf{H} = (\mathbf{h}(\mathbf{x}^1), \dots, \mathbf{h}(\mathbf{x}^K)), \quad \overline{\nabla_{\mathbf{x}}\mathbf{h}} = \frac{1}{K} \sum_{k=1}^K \nabla_{\mathbf{x}^k}\mathbf{h}(\mathbf{x}^k),$$

$$\mathbf{G} := (\nabla_{\mathbf{x}^1} \log q(\mathbf{x}^1), \dots, \nabla_{\mathbf{x}^K} \log q(\mathbf{x}^K))^\top,$$

$$\text{Then } -\frac{1}{K}\mathbf{H}\mathbf{G} + \text{err} = \overline{\nabla_{\mathbf{x}}\mathbf{h}}.$$

Stein gradient estimator (kernel based)

Main idea: invert Stein's identity:

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

Matrix form: $-\frac{1}{K}\mathbf{H}\mathbf{G} + \text{err} = \overline{\nabla_{\mathbf{x}}\mathbf{h}}$.

3. Now solve a ridge regression problem:

$$\hat{\mathbf{G}}_V^{\text{Stein}} := \arg \min_{\hat{\mathbf{G}} \in \mathbb{R}^{K \times d}} \|\overline{\nabla_{\mathbf{x}}\mathbf{h}} + \frac{1}{K}\mathbf{H}\hat{\mathbf{G}}\|_F^2 + \frac{\eta}{K^2}\|\hat{\mathbf{G}}\|_F^2,$$

Stein gradient estimator (kernel based)

Main idea: invert Stein's identity:

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

Matrix form: $-\frac{1}{K}\mathbf{H}\mathbf{G} + \text{err} = \overline{\nabla_{\mathbf{x}}\mathbf{h}}$.

3. Now solve a ridge regression problem:

$$\hat{\mathbf{G}}_V^{\text{Stein}} := \arg \min_{\hat{\mathbf{G}} \in \mathbb{R}^{K \times d}} \|\overline{\nabla_{\mathbf{x}}\mathbf{h}} + \frac{1}{K}\mathbf{H}\hat{\mathbf{G}}\|_F^2 + \frac{\eta}{K^2}\|\hat{\mathbf{G}}\|_F^2,$$

Analytic solution: $\hat{\mathbf{G}}_V^{\text{Stein}} = -(\mathbf{K} + \eta\mathbf{I})^{-1}\langle \nabla, \mathbf{K} \rangle,$

with

$$\mathbf{K} := \mathbf{H}^{\top}\mathbf{H}, \quad \mathbf{K}_{ij} = \mathcal{K}(\mathbf{x}^i, \mathbf{x}^j) := \mathbf{h}(\mathbf{x}^i)^{\top}\mathbf{h}(\mathbf{x}^j),$$

$$\langle \nabla, \mathbf{K} \rangle := K\mathbf{H}^{\top}\overline{\nabla_{\mathbf{x}}\mathbf{h}}, \quad \langle \nabla, \mathbf{K} \rangle_{ij} = \sum_{k=1}^K \nabla_{\mathbf{x}_j^k} \mathcal{K}(\mathbf{x}^i, \mathbf{x}^k).$$

Gradient estimators: comparisons

Comparing KDE plugin gradient estimator and Stein gradient estimator:
for translation invariant kernels:

$$\hat{\mathbf{G}}^{\text{KDE}} = -\text{diag}(\mathbf{K}\mathbf{1})^{-1} \langle \nabla, \mathbf{K} \rangle$$

$$\hat{\mathbf{G}}_{\text{V}}^{\text{Stein}} = -(\mathbf{K} + \eta \mathbf{I})^{-1} \langle \nabla, \mathbf{K} \rangle$$

When approximating $\nabla_{\mathbf{x}^k} \log q(\mathbf{x}^k)$:

- KDE: only use $\mathcal{K}(\mathbf{x}^j, \mathbf{x}^k)$
- Stein: use all $\mathcal{K}(\mathbf{x}^j, \mathbf{x}^i)$ even for those $i \neq k$

more sample efficient!

Gradient estimators: comparisons

Compare to the score matching gradient estimator:

Score matching

- Min. expected l_2 error
(a stronger divergence)
- Parametric approx.
(introduce approx. error)
- Repeated derivations for
different kernels

Stein

- Min. KSD
(a weaker divergence)
- Non-parametric approx.
(no approx. error)
- Ubiquitous solution for
any kernel

KSD: Kernelised Stein discrepancy (Liu et al. 2016; Chwialkowski et al. 2016)

Example: meta-learning for approximate inference

- learn an approx. posterior sampler for NN weights

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \zeta \Delta_{\phi}(\boldsymbol{\theta}_t, \nabla_t) + \boldsymbol{\sigma}_{\phi}(\boldsymbol{\theta}_t, \nabla_t) \odot \boldsymbol{\epsilon}_t, \quad \boldsymbol{\epsilon}_t \sim \mathcal{N}(\boldsymbol{\epsilon}; \mathbf{0}, \mathbf{I}),$$

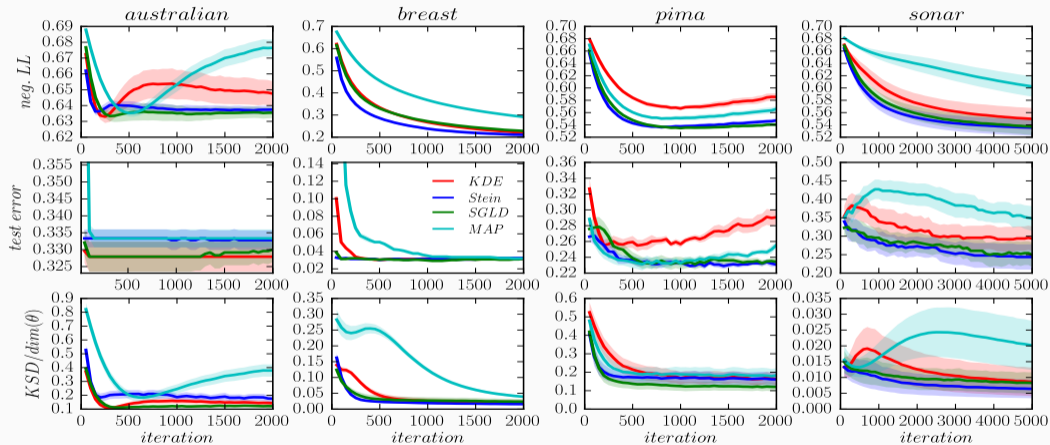
$$\nabla_t = \nabla_{\boldsymbol{\theta}_t} \left[\frac{N}{M} \sum_{m=1}^M \log p(y_m | \mathbf{x}_m, \boldsymbol{\theta}_t) + \log p_0(\boldsymbol{\theta}_t) \right].$$

- coordinates of $\Delta_{\phi}(\boldsymbol{\theta}_t, \nabla_t)$ and $\boldsymbol{\sigma}_{\phi}(\boldsymbol{\theta}_t, \nabla_t)$ are parameterised by MLPs
- training objective: an MC approximation of

$$\sum_t \mathcal{L}_{VI}(q_t), \quad q_t \text{ is the marginal distribution of } \boldsymbol{\theta}_t$$

- see whether it generalises to diff. architectures and datasets:
 - train: 1-hidden-layer BNN with 20 hidden units + ReLU, on *crabs* dataset
 - test: 1-hidden-layer BNN with 50 hidden units + sigmoid, on other datasets

Example: meta-learning for approximate inference



What we covered today:

- Is density evaluation really necessary for inference tasks?
- Fitting implicit approx. posterior
by **approximating variational lower-bound's gradients**
- Designing implicit posterior approximations: big challenge

Thank you!

(BDL workshop tomorrow: training implicit generative models)