## Wild Approximate Inference: Why and How

Yingzhen Li

Machine Learning Group, University of Cambridge Joint work with Rich Turner (Cambridge) and Qiang Liu (Dartmouth  $\rightarrow$  UT Austin)

- 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."

- 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?

• In Bayesian modelling, we specify the model:

prior: p(z), likelihood: p(x|z)

• We want to compute the exact posterior

exact posterior: 
$$p(z|x) = rac{p(z)p(x|z)}{p(x)}$$

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

- In most of the time we cannot compute p(z|x) efficiently
- Approximate inference: find q(z|x) in some family Q such that  $q(z|x) \approx p(z|x)$
- and in inference time, approximate

$$\mathbb{E}_{p(\boldsymbol{z}|\boldsymbol{x})}\left[F(\boldsymbol{z})\right] \approx \mathbb{E}_{q(\boldsymbol{z}|\boldsymbol{x})}\left[F(\boldsymbol{z})\right]$$

- In most of the time we cannot compute p(z|x) efficiently
- Approximate inference: find q(z|x) in some family Q such that  $q(z|x) \approx p(z|x)$
- and in inference time, approximate

$$\mathbb{E}_{p(z|x)}[F(z)] \approx \mathbb{E}_{q(z|x)}[F(z)]$$

## Tractability: fast computation of the RHS term

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

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

$$q^*(\boldsymbol{z}|\boldsymbol{x}) = \argmin_{q \in \mathcal{Q}} \operatorname{KL}[q(\boldsymbol{z}|\boldsymbol{x})||p(\boldsymbol{z}|\boldsymbol{x})] = \arg_{q \in \mathcal{Q}} \operatorname{KL}[q(\boldsymbol{z}|\boldsymbol{x})] + \mathbb{H}[q(\boldsymbol{z}|\boldsymbol{x})]$$

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

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

• With Monte Carlo approximation methods, inference is done by

$$\mathbb{E}_{p(\boldsymbol{z}|\boldsymbol{x})}\left[F(\boldsymbol{z})
ight] pprox rac{1}{K}\sum_{k=1}^{K}F(\boldsymbol{z}^k), \quad \boldsymbol{z}^k \sim q(\boldsymbol{z}|\boldsymbol{x})$$

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

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

$$q^*(oldsymbol{z}|oldsymbol{x}) = rgmin_{q\in\mathcal{Q}} \operatorname{KL}[q(oldsymbol{z}|oldsymbol{x})||p(oldsymbol{z}|oldsymbol{x})] = rgmax_{q\in\mathcal{Q}} \operatorname{E}_q[\log p(oldsymbol{z},oldsymbol{x})] + \mathbb{H}[q(oldsymbol{z}|oldsymbol{x})]$$

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

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

• With Monte Carlo approximation methods, inference is done by

$$\mathbb{E}_{p(\boldsymbol{z}|\boldsymbol{x})}\left[F(\boldsymbol{z})
ight] pprox rac{1}{K}\sum_{k=1}^{K}F(\boldsymbol{z}^k), \quad \boldsymbol{z}^k \sim q(\boldsymbol{z}|\boldsymbol{x})$$

Tractability requirement: fast sampling

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

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

$$q^*(oldsymbol{z}|oldsymbol{x}) = rgmin_{q\in\mathcal{Q}} \operatorname{KL}[q(oldsymbol{z}|oldsymbol{x})||p(oldsymbol{z}|oldsymbol{x})] = rgmax_{q\in\mathcal{Q}} \operatorname{E}_q[\log p(oldsymbol{z},oldsymbol{x})] + \mathbb{H}[q(oldsymbol{z}|oldsymbol{x})]$$

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

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

• With Monte Carlo approximation methods, inference is done by

$$\mathbb{E}_{p(\boldsymbol{z}|\boldsymbol{x})}\left[F(\boldsymbol{z})
ight] pprox rac{1}{K}\sum_{k=1}^{K}F(\boldsymbol{z}^k), \quad \boldsymbol{z}^k \sim q(\boldsymbol{z}|\boldsymbol{x})$$

Tractability requirement: fast sampling and fast density (ratio) evaluation (only for optimisation!) 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 do not require density evaluation!

Research for wild approximate inference:

Can we design efficient approximate inference algorithms that enables fast inference, without adding more constraints to q? Research for wild approximate inference:

# Can we design efficient approximate inference algorithms that enables fast inference, without adding more constraints 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
- Fast inference
- Easy to amortise (memory efficient)

#### MCMC

- Just need to do sampling
- Very accurate
- Need large T thus slow (MCMC)
- Not re-usable when *p* is updated

We want to have the best from both worlds!

#### Meta learning for approximate inference:

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

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)
- Gradient approximations Huszár (2017); Li and Turner (2017)
- Amortising deterministic/stochastic dynamics Wang and Liu (2016); Li, Turner and Liu (2017); Chen et al. (2017); Pu et al. (2017)

See Chapter 5 of my draft thesis for a discussion Also see Liu and Lee (2017), Titsias (2017)

## Variational lower-bound approximation via density ratio estimation

Recall the variational lower-bound

$$\mathcal{L}_{\mathsf{VI}}(q) = \mathbb{E}_{q}\left[\lograc{p(oldsymbol{x}|oldsymbol{z})p(oldsymbol{z})}{q(oldsymbol{z}|oldsymbol{x})}
ight]$$

- let's assume  $p(\mathbf{x}|\mathbf{z})$  is tractable
- then it only remains to estimate the density ratio  $\frac{p(z)}{q(z|x)}$
- classifier-based GAN (Goodfellow et al. 2014): the optimal discriminator satisfies

$$D^*(\boldsymbol{z}, \boldsymbol{x}) = \left[1 + \exp\left[-\log rac{p(\boldsymbol{z})}{q(\boldsymbol{z}|\boldsymbol{x})}
ight]
ight]^{-1}$$

Therefore,

$$\mathcal{L}_{\mathsf{VI}}(q) pprox \mathbb{E}_q\left[\log p(oldsymbol{x}|oldsymbol{z}) + \log rac{D(oldsymbol{z},oldsymbol{x})}{1 - D(oldsymbol{z},oldsymbol{x})}
ight].$$

In practice need to do "adaptive contrast" (Mescheder et al. 2017) Can extend this to IS/SMC (NIPS workshop) Minimising Stein's discrepancy

$$\mathcal{S}^2(p(oldsymbol{z}|oldsymbol{x}),q(oldsymbol{z}|oldsymbol{x})) = \left(\sup_{h\in\mathcal{H}}\mathbb{E}_q\left[
abla_z\log p(oldsymbol{z},oldsymbol{x})h(oldsymbol{z}) + \langle 
abla,h(oldsymbol{z})
ight)^2$$

- this is nice: only need samples from q and evaluate the joint distribution p(z, x)
- OPVI (Ranganath et al. 2016):  $\mathcal{H}$  is a set of neural network functions
- Wild VI via KSD minimisation (Liu and Feng 2016):  ${\cal H}$  is the unit ball in an RKHS

## **Gradient Approximations**

## Alternative idea: approximate the gradient

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

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

During optimisation we only care about the gradients!

The gradient of the variational lower-bound:

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

## Alternative idea: approximate the gradient

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

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

During optimisation we only care about the gradients!

The gradient of the variational lower-bound:

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

The gradient of the entropy term:

$$\nabla_{\phi} \mathbb{H}[q(\boldsymbol{z}|\boldsymbol{x})] = -\mathbb{E}_{\pi} \left[ \nabla_{\boldsymbol{f}} \log q(\boldsymbol{f}_{\phi}(\boldsymbol{\epsilon}, \boldsymbol{x})|\boldsymbol{x})^{\mathsf{T}} \nabla_{\phi} \boldsymbol{f}_{\phi}(\boldsymbol{\epsilon}, \boldsymbol{x}) \right] - \overline{\mathbb{E}_{q} \left[ \nabla_{\phi} \log q_{\phi}(\boldsymbol{z}|\boldsymbol{x}) \right]}$$
  
this term is 0

It remains to approximate  $\nabla_z \log q(z|x)!$ (in a cheap way) KDE plug-in gradient estimator for  $abla_{m{x}}\log q(m{x})$  for  $m{x}\in\mathbb{R}^d$ :

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

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

• then approximate  $abla_{m{x}}\log q(m{x}) pprox 
abla_{m{x}}\log \hat{q}(m{x})$ 



### Gradient estimators (kernel based)

Score matching gradient estimators: find  $\hat{g}(x)$  to minimise the  $\ell_2$  error

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

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

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

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

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

and find the best  ${m a}=(a_1,...,a_K)$  by minimising the  $\ell_2$  error.

Define h(x): a (column vector) test function satisfying the boundary condition

 $\lim_{x\to\infty}q(x)\boldsymbol{h}(x)=\boldsymbol{0}.$ 

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

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

Define h(x): a (column vector) test function satisfying the boundary condition

 $\lim_{x\to\infty}q(x)\boldsymbol{h}(x)=\boldsymbol{0}.$ 

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

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

Invert Stein's identity to obtain  $\nabla_x \log q(x)$ !

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. MC approximation to Stein's identity:

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

## Stein gradient estimator (kernel based)

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. MC approximation to Stein's identity:

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

2. Rewrite the MC equations in matrix forms: denoting

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

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:  $-\frac{1}{\kappa}\mathbf{H}\mathbf{G} + \mathrm{err} = \overline{\nabla_{\mathbf{x}}\mathbf{h}}.$ 

3. Now solve a ridge regression problem:

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

## Stein gradient estimator (kernel based)

Main idea: invert Stein's identity:

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

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

3. Now solve a ridge regression problem:

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

Analytic solution:  $\hat{\mathbf{G}}_V^{ ext{Stein}} = -(\mathbf{K} + \eta \mathbf{I})^{-1} \langle 
abla, \mathbf{K} 
angle$ , with

$$\begin{array}{ll} \mathsf{K} := \mathsf{H}^{\mathsf{T}}\mathsf{H}, \quad \mathsf{K}_{ij} = \mathcal{K}(\boldsymbol{x}^{i}, \boldsymbol{x}^{j}) := \boldsymbol{h}(\boldsymbol{x}^{i})^{\mathsf{T}}\boldsymbol{h}(\boldsymbol{x}^{j}), \\ \langle \nabla, \mathsf{K} \rangle := \mathcal{K}\mathsf{H}^{\mathsf{T}}\overline{\nabla_{\boldsymbol{x}}\boldsymbol{h}}, \quad \langle \nabla, \mathsf{K} \rangle_{ij} = \sum_{k=1}^{K} \nabla_{\boldsymbol{x}_{j}^{k}} \mathcal{K}(\boldsymbol{x}^{i}, \boldsymbol{x}^{k}). \end{array}$$

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

$$\hat{\mathsf{G}}^{\mathsf{KDE}} = -\mathsf{diag}(\mathsf{K1})^{-1} \langle 
abla, \mathsf{K} 
angle$$

$$\hat{\mathbf{G}}_V^{ ext{Stein}} = -(\mathbf{K} + \eta \mathbf{I})^{-1} \langle 
abla, \mathbf{K} 
angle$$

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!

Compare to the score matching gradient estimator:

Score matching

- Min. expected l<sub>2</sub> 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

Complexity: both need matrix inversion,  $\mathcal{O}(K^3 + K^2 d)$ 

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

Use Stein gradient estimator at unseen locations y:

- non-parametric predictive model:
  - compute  $\hat{\mathbf{G}}_V^{ ext{Stein}}$  on  $\{\mathbf{y}\} \cup \{\mathbf{x}^k\}_{k=1}^K$
  - retrieve the row corresponding to the approximation of  $\log q(y)$
- parametric predictive model:
  - define  $\hat{g}(x) = \sum_{k=1}^{K} a_k \nabla_x \mathcal{K}(x, x^k), \quad x^k \sim q(x)$
  - find the best  $\boldsymbol{a} = (a_1, ..., a_K)$  by minimising KSD
  - compute  $\hat{g}(y)$  using the fitted *a* and the stored samples  $\{x^k\}_{k=1}^K$

### Example: meta-learning for approximate inference

• learn an approx. posterior sampler for NN weights

$$\begin{aligned} \boldsymbol{\theta}_{t+1} &= \boldsymbol{\theta}_t + \zeta \Delta_{\boldsymbol{\phi}}(\boldsymbol{\theta}_t, \nabla_t) + \boldsymbol{\sigma}_{\boldsymbol{\phi}}(\boldsymbol{\theta}_t, \nabla_t) \odot \boldsymbol{\epsilon}_t, \quad \boldsymbol{\epsilon}_t \sim \mathcal{N}(\boldsymbol{\epsilon}; \boldsymbol{0}, \boldsymbol{I}) \\ \nabla_t &= \nabla_{\boldsymbol{\theta}_t} \left[ \frac{N}{M} \sum_{m=1}^M \log p(y_m | \boldsymbol{x}_m, \boldsymbol{\theta}_t) + \log p_0(\boldsymbol{\theta}_t) \right]. \end{aligned}$$

- coordinates of  $\Delta_{\phi}( heta_t, 
  abla_t)$  and  $\sigma_{\phi}( heta_t, 
  abla_t)$  are parameterised by MLPs
- training objective: an MC approximation of

$$\sum_t \mathcal{L}_{\mathsf{VI}}(q_t), \quad q_t$$
 is the marginal distribution of  $oldsymbol{ heta}$ 

- 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



• Addressing mode collapse: train your generator using entropy regularisation:

min  $\mathcal{L}_{gen}(p_{gen}) - \mathbb{H}[p_{gen}]$ 

- $\mathcal{L}_{gen}(p_{gen})$  is the generator loss of your favourite GAN method
- Again the gradient of  $\mathbb{H}[p_{\mathsf{gen}}]$  is approximated by the gradient estimators
- We pick BEGAN because it has severe mode collapse problem

## Example: entropy regularised GANs



• training a generative model  $q_{\theta}(x)$  by minimising

 $\min_{\theta} \mathrm{KL}[q_{\theta}(\mathbf{x}) || p_{\mathcal{D}}(\mathbf{x})]$ 

• the gradient of KL: say  $\pmb{x}\sim \pmb{q}\Leftrightarrow \pmb{\epsilon}\sim \pi(\pmb{\epsilon}), \pmb{x}=\pmb{f}_{\pmb{ heta}}(\pmb{\epsilon})$ 

$$\nabla_{\boldsymbol{\theta}} \mathrm{KL} = \mathbb{E}_{\pi}[(\nabla_{\boldsymbol{x}} \log q_{\boldsymbol{\theta}}(\boldsymbol{x}) - \nabla_{\boldsymbol{x}} \log p_{\mathcal{D}}(\boldsymbol{x}))|_{\boldsymbol{x} = \boldsymbol{f}_{\boldsymbol{\theta}}(\boldsymbol{\epsilon})} \nabla_{\boldsymbol{\theta}} \boldsymbol{f}]$$

• Use the gradient estimator for both  $\nabla_x \log q_\theta(x)$  and  $\nabla_x \log p_D(x)!$ 

### Example: training implicit generative models



## Amortised MCMC

## Recap: from Variational EM to VAE

## Variational EM



(For simplicity we omit the model parameter  $\theta$  but it would be trained by approx. MLE)

## Recap: from Variational EM to VAE



Amortised inference: memory efficient & no need to run VI optimisation in test time!

#### MCMC-EM



For every  $x_n$ , need to simulate MCMC for T >> 0 steps (slow!)

## New: from MCMC-EM to amortised MCMC



Need to store all samples from the previous iteration, memory cost O(NKD). For a new datapoint, still need to run MCMC with T >> 0 starting from  $p_0$  (slow!)

## New: from MCMC-EM to amortised MCMC



This method essentially cares about  $q_T$  only, so no need for  $q(z|x) \approx p(z|x)$ . In test time still need to run MCMC to obtain samples from  $q_T$  (slow!)

## New: from MCMC-EM to amortised MCMC



Distillation happens at the same time during training (thus also improving  $q_T$ ), and now  $q(z|x) \approx p(z|x)$  – no need to run MCMC in test time!

## Understanding "distillation during training"

I've generated some samples. Do they approximate the target distribution well?

You can do better. Returned you some MCMC samples improved upon yours.

OK thanks I'll try to generate these improved samples next time.





Sampling-based methods, e.g. Markov chain Monte Carlo:

- Define a transition kernel  $\mathcal{K}(\boldsymbol{z}_t | \boldsymbol{z}_{t-1})$  that leaves  $p(\boldsymbol{z} | \boldsymbol{x})$  invariant
- Pick an initial distribution  $q_0(z_0|x)$ , and run *T*-step transitions

$$oldsymbol{z}_T^k \sim q_T(oldsymbol{z}_T | oldsymbol{x}) = \int q_0(oldsymbol{z}_0 | oldsymbol{x}) \prod_{t=1}^T \mathcal{K}(oldsymbol{z}_T | oldsymbol{z}_{t-1}) doldsymbol{z}_{0:T-1}, \quad k=1,2,...,K$$

- We rely on two assumptions:
  - $\mathcal{K}(z_t|z_{t-1})$  has a unique stationary distribution p(z|x)
  - when  $\mathcal{T} 
    ightarrow +\infty$  the Markov chain converges to the stationary
- Then for all T > 0,  $q_0(z|x) = q_T(z|x)$  iff.  $q_0(z|x) = p(z|x)$

Three main ingredients of the framework:

- "Student": an architecture of q(z|x)
- "Teacher": a transition kernel K invariant to p(z|x), which produces q<sub>T</sub>(z|x)
- "Feedback": an update rule to match q(z|x) to q<sub>T</sub>(z|x) (say divergence minimisation)

After looping,  $q(\boldsymbol{z}|\boldsymbol{x}) \approx p(\boldsymbol{z}|\boldsymbol{x})$ , and use  $q(\boldsymbol{z}|\boldsymbol{x})$  for inference!



• TD-learning and Q-learning:

also tries to update the current V- or Q-value to using T-step roll-out

- Can borrow RL ideas for convergence proof: TD/Q-learning can show convergence when using linear approximations
- using NNs as  $q_0 \leftrightarrow$  using DQNs
- Can borrow deep RL tricks to design distillation rules!

#### We want q to be implicit!

(i.e. can sample from q but cannot evaluate density)

idea: match samples  $\{ \pmb{z}_0^k \} \sim q$  to samples  $\{ \pmb{z}_T^k \} \sim q_T !$ 

- We tested the original GAN idea  $^1$
- In general, any GAN-like technique is applicable!

<sup>&</sup>lt;sup>1</sup>Goodfellow et al. Generative Adversarial Networks. NIPS 2014.

Generative model: a small convolutional decoder

- VAE training:
  - Gaussian encoder: a symmetric flip of the decoder
- Amortised MCMC (AMC) training:
  - CNN-G encoder:  $z = MLP([CNN(x), \epsilon]), \epsilon \sim \mathcal{N}(\mathbf{0}, I)$
  - CNN-B encoder:  $z = MLP([CNN(x) \odot \epsilon]), \epsilon \sim Bern(0.5)$
  - MCMC: Langevin Dynamics w/out rejection

0/23+56789 0/23456789 0/23456789 0/23456789 0/23456789 0/23456789 0/23456789 0/23456789 0/23456789 0/23456789 0/23456789	2 0 0 0 4 7 4 9 4 6 7 5 7 8 9 5 9 0 0 7 3 4 8 4 5 4 7 9 7 4 9 5 4 5 4 7 9 7 7 7 2 7 2 7 2 9 7 3 4 9 5 4 5 4 7 9 7 7 7 2 3 9 2 7 2 9 9 3 4 0 5 4 7 8 2 7 2 9 2 7 2 9 9 3 4 0 5 4 7 8 2 7 2 9 2 7 2 9 9 3 4 0 5 4 7 8 2 4 9 7 4 7 6 4 3 6 1 6 3 9 4 7 0 9 7 7 6 4 6 Gaussian encoder + VAE
7 9 9 5 5 7 6 9 6 5	0 0 4 7 3 7 3 7 3 7 4 1
3 9 7 9 1 9 1 4 2 6 2 7	8 9 0 5 4 4 2 3 9 4 3
7 9 9 9 4 4 0 5 3 2	7 1 9 1 2 3 4 2 3 5
2 / 4 9 1 0 3 9 6 5	7 5 3 6 7 4 4 2 1 4
3 9 3 2 8 1 5 1	4 9 6 9 2 4 8 3 4 2
3 9 1 7 9 8 1 7 1 9 6	3 3 0 7 4 6 3
4 7 7 9 8 1 7 1 9 6	1 4 5 5 3 0 1 4 4 6
6 9 5 2 1 5 7 6 7 4	7 8 6 9 3 6 9 8 4
7 7 1 9 4 5 5 6 6	4 3 7 4 // 2 1 4 5
6 9 2 4 5 5 6 7	9 6 0 2 6 9 7 7 2
CNN-G + AMC	CNN-B + AMC

Table 1: Average Test LL and effective sample size (ESS).  $\eta$  as the stepsize for Langevin dynamics.

Encoder	Method	IW-LL	IW-ESS	HAIS-LL	HAIS-ESS
Gaussian	VAE	-81.31	104.11	-80.64	91.59
	MCMC-VI, $T=$ 5, $\eta=$ 0.2	-90.06	110.58	-89.79	85.63
	AMC, $T=$ 5, $\eta=$ 0.2	-90.71	49.02	-89.64	87.93
CNN-G	AMC, $T=$ 5, $\eta=$ 0.2	-90.84	31.60	-89.35	87.49
	AMC, $T=$ 50, $\eta=$ 0.02	-83.30	6.84	-78.23	77.78
	AVB	-94.97	11.30	-85.92	57.21
CNN-B	AMC, $T=$ 5, $\eta=$ 0.2	-90.75	34.17	-89.42	88.10
	AMC, $T=$ 50, $\eta=$ 0.02	-83.62	8.88	-80.03	80.71
	AVB	-89.47	8.98	-82.66	76.90
N/A	persistent MCMC, $T=50,~\eta=0.02$	-84.43	9.14	-78.88	77.29

• HAIS seems to be more reliable (K = 100), compared to importance sampling (IS, K = 5000)

• The best case (CNN-G) is better than persistent MCMC

## Missing data imputation

- Given an image  $\mathbf{x} = [\mathbf{x}_o, \mathbf{x}_m]$  with missing values, repeat the following for T steps:
  - sample  $z \sim q(z|x_o, x_m)$
  - sample  $x^* \sim p(x|z)$  and set  $x_m \leftarrow x_m^*$

**Table 2:** Label entropy on nearest neighbours. The  $l_1$  distance is divided by the number of pixels.

Dataset	VAE	CNN-G	CNN-B
Entropy	$0.411{\pm}0.0389$	$0.701{\pm}0.0476$	$0.933{\pm}0.0491$
<i>l</i> <sub>1</sub> -norm	$0.061{\pm}0.0002$	$0.059 {\pm} 0.0001$	$0.064{\pm}0.0002$



#### We want q to be implicit!

If you don't like discriminators/minimax optimisation: do moment matching!

- Simply match the first *M* moments (not minimising a divergence)
  - How to choose *M* and how to reweigh the importance of them?
- We propose energy matching:

$$\min_{\phi} \left| \mathbb{E}_{q_{\mathcal{T}}} \left[ \log p(\boldsymbol{x}, \boldsymbol{z}_{\mathcal{T}}) \right] - \mathbb{E}_{q_{\phi}} \left[ \log p(\boldsymbol{x}, \boldsymbol{z}) \right] \right|^2$$

- Let the log-joint distribution tell you which moments you want to match
- Related to Contrastive Divergence (CD-T), works well with small T

Table 3: BNN classification experiment results. MALA = Langevin dynamics with rejection

	Average Test Log-likelihood			Average Test Error		
Dataset	VI+Gaussian	AMC	MALA	VI+Gaussian	AMC	MALA
australian	$-0.633 \pm 0.008$	$-0.666 \pm 0.015$	$-0.636 \pm 0.010$	$0.315 {\pm} 0.014$	$0.360 \pm 0.011$	$0.344 \pm 0.013$
breast	$-0.096 \pm 0.010$	$-0.091 \pm 0.008$	$-0.094 \pm 0.010$	$0.029 {\pm} 0.003$	$0.030 \pm 0.004$	$0.037 \pm 0.004$
colon	$-0.799 \pm 0.246$	$-0.491 \pm 0.104$	-0.420±0.027	$0.125 \pm 0.023$	$0.167 \pm 0.031$	$0.167 \pm 0.026$
crabs	$-0.221 \pm 0.012$	$-0.115 \pm 0.011$	$-0.179 \pm 0.010$	$0.070 \pm 0.013$	$0.040 \pm 0.010$	$0.035 \pm 0.010$
ionosphere	$-0.241 \pm 0.019$	$-0.230 \pm 0.031$	$-0.179 \pm 0.013$	$0.099 \pm 0.011$	$0.077 \pm 0.013$	$0.064 \pm 0.010$
pima	$-0.503 \pm 0.010$	$-0.506 \pm 0.013$	$-0.498 \pm 0.012$	$0.262 \pm 0.008$	$0.245 \pm 0.008$	$0.247 \pm 0.008$
sonar	$-0.389 \pm 0.025$	$-0.347 \pm 0.030$	$-0.366 \pm 0.021$	$0.179 \pm 0.014$	$0.150 \pm 0.016$	$0.171 \pm 0.020$



What we covered today:

- Is density evaluation really necessary for inference tasks?
- Fitting implicit approx. posterior:
  - objective approximation via density ratio estimation
  - minimising Stein's discrepancy
  - gradient approximations
  - amortise MCMC/SVGD
- Designing implicit posterior approximations: big challenge

## Thank you!