

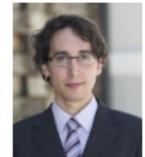


Gradient Estimation for Implicit Models with Stein's Method

Yingzhen Li

Microsoft Research Cambridge

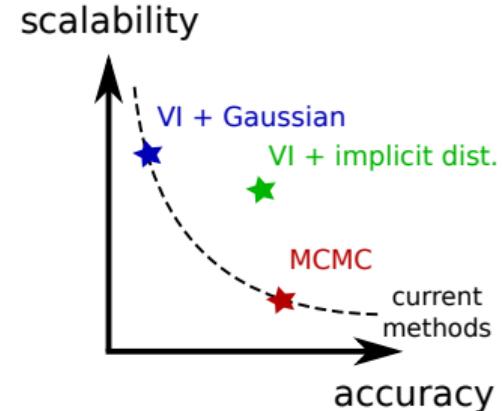
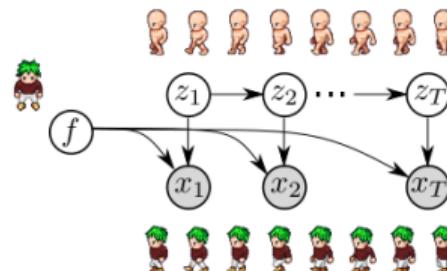
Joint work with Rich Turner, Wenbo Gong,
and José Miguel Hernández-Lobato



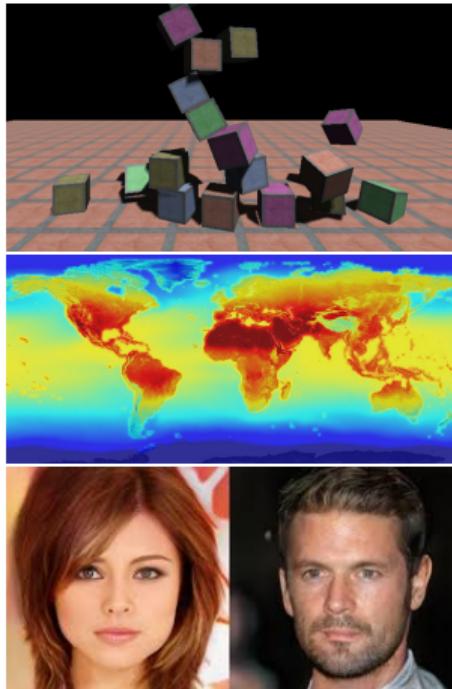
A little about my research...



Bayesian Deep Learning



Examples for implicit (generative) models



Implicit distributions:

- + easy to sample:
 $\mathbf{z} \sim q(\mathbf{z}|\mathbf{x}) \Leftrightarrow \epsilon \sim \pi(\epsilon), \mathbf{z} = f(\epsilon, \mathbf{x})$
- + super flexible

Examples for implicit (generative) models

- Bayesian inference goal: compute $\mathbb{E}_{p(z|x)} [F(z)]$
- Approximate inference: find $q(z|x)$ in some family \mathcal{Q} such that $q(z|x) \approx p(z|x)$
- At inference time: Monte Carlo integration:

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

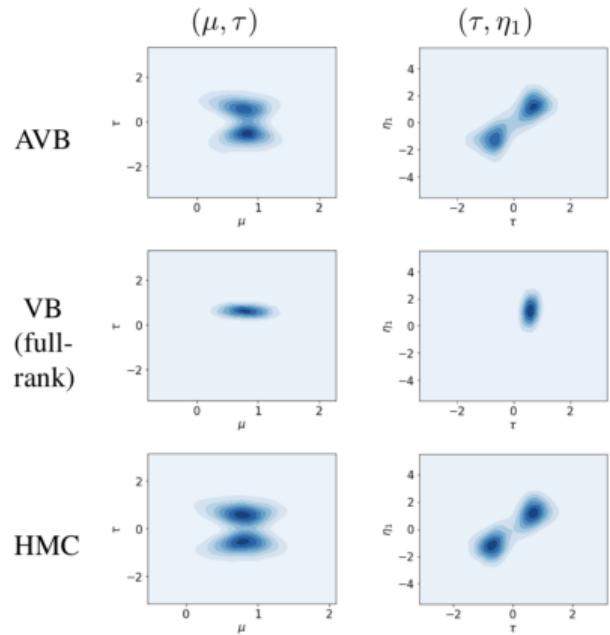
Examples for implicit (generative) models

- Bayesian inference goal: compute $\mathbb{E}_{p(z|x)} [F(z)]$
- Approximate inference: find $q(z|x)$ in some family \mathcal{Q} such that $q(z|x) \approx p(z|x)$
- At inference time: Monte Carlo integration:

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

Tractability requirement: fast sampling from q
(no need for point-wise density evaluation)

Examples for implicit (generative) models

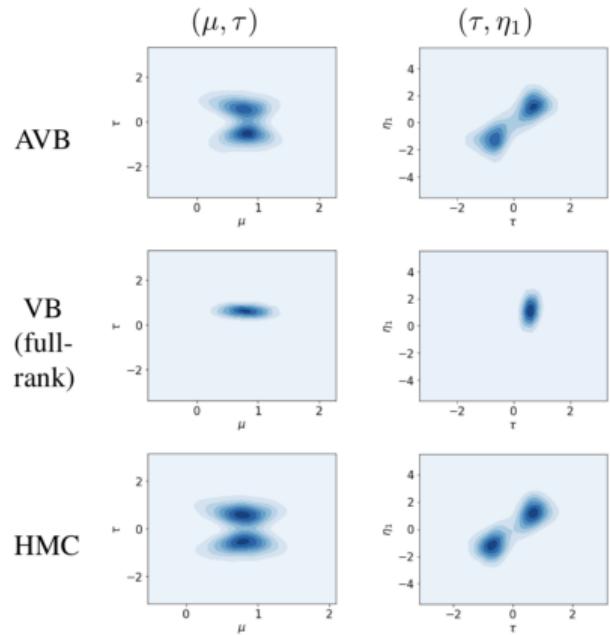


Implicit distributions:

- + easy to sample:
 $\mathbf{z} \sim q(\mathbf{z}|\mathbf{x}) \Leftrightarrow \epsilon \sim \pi(\epsilon), \mathbf{z} = \mathbf{f}(\epsilon, \mathbf{x})$
- + super flexible
- + better approximate posterior

Fig. source: Mescheder et al. (2017)

Examples for implicit (generative) models



Implicit distributions:

- + easy to sample:
 $\mathbf{z} \sim q(\mathbf{z}|\mathbf{x}) \Leftrightarrow \epsilon \sim \pi(\epsilon), \mathbf{z} = \mathbf{f}(\epsilon, \mathbf{x})$
- + super flexible
- + better approximate posterior
- hard to evaluate density,
need some tricks for training

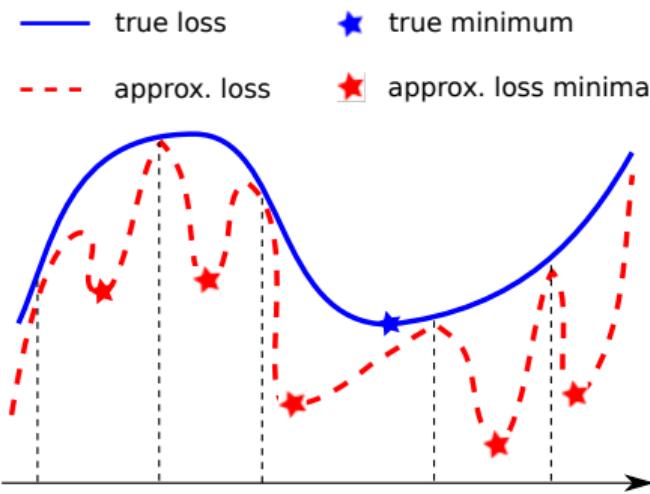
Fig. source: Mescheder et al. (2017)

Loss approximation vs gradient approximation

To train the implicit generative model $p_\phi(x)$:

E.g. the generative adversarial network (GAN) method (Goodfellow et al. 2014):

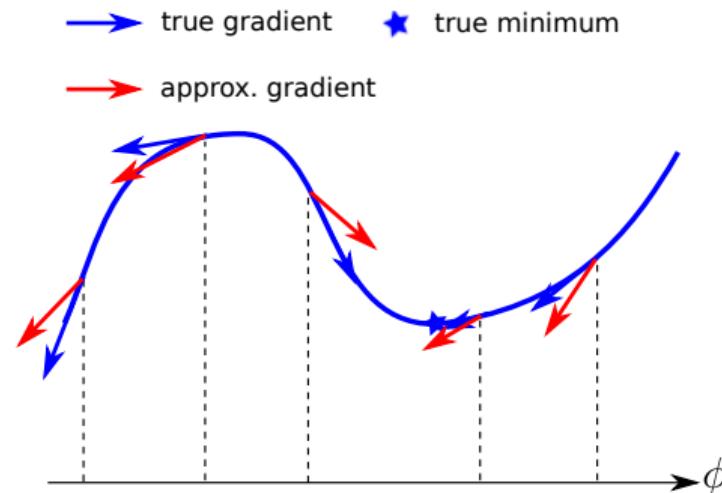
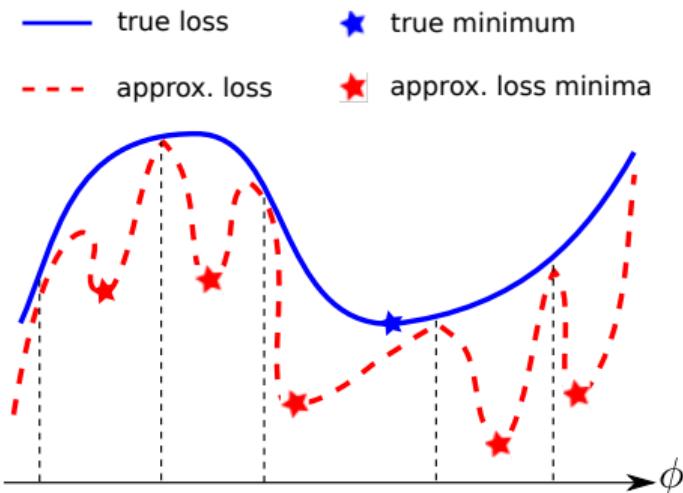
$$\min_{\phi} \text{JS}[p_D || p_\phi] = \min_{\theta} \max_D \mathbb{E}_{p_D} [\log D(x)] + \mathbb{E}_{p_\phi} [\log(1 - D(x))]$$



Loss approximation vs gradient approximation

Often we use gradient-based optimisation methods to train machine learning models.

... which only require evaluating the gradient, rather than the loss function!



Gradient approximation for VI

Variational inference with q distribution parameterised by ϕ :

$$\phi^* = \arg \min_{\phi} \text{KL}[q_{\phi}(z|x) || p(z|x)] = \arg \max_{\phi} \mathcal{L}_{\text{VI}}(q_{\phi})$$

$$\begin{aligned}\mathcal{L}_{\text{VI}}(q_{\phi}) &= \log p(x) - \text{KL}[q_{\phi}(z|x) || p(z|x)] \\ &= \log p(x) - \mathbb{E}_q \left[\log \frac{q_{\phi}(z|x)}{p(z|x)} \right] \\ &= \mathbb{E}_q \left[\log \frac{p(x, z)}{q_{\phi}(z|x)} \right] \\ &= \mathbb{E}_q [\log p(x, z)] + \mathbb{H}[q_{\phi}(z|x)]\end{aligned}$$

$\mathcal{L}_{\text{VI}}(q_{\phi})$ is also called the variational lower-bound

Gradient approximation for VI

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

$$\begin{aligned}\mathcal{L}_{\text{VI}}(q_\phi) &= \mathbb{E}_q [\log p(\mathbf{x}, \mathbf{z})] + \mathbb{H}[q_\phi(\mathbf{z}|\mathbf{x})] \\ &= \mathbb{E}_\pi [\log p(\mathbf{x}, \mathbf{f}_\phi(\epsilon, \mathbf{x}))] + \mathbb{H}[q_\phi(\mathbf{z}|\mathbf{x})] \quad // \text{ reparam. trick}\end{aligned}$$

Gradient approximation for VI

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

$$\begin{aligned}\mathcal{L}_{\text{VI}}(q_\phi) &= \mathbb{E}_q [\log p(\mathbf{x}, \mathbf{z})] + \mathbb{H}[q_\phi(\mathbf{z}|\mathbf{x})] \\ &= \mathbb{E}_\pi [\log p(\mathbf{x}, \mathbf{f}_\phi(\epsilon, \mathbf{x}))] + \mathbb{H}[q_\phi(\mathbf{z}|\mathbf{x})] \quad // \text{ reparam. trick}\end{aligned}$$

If you use gradient descent for optimisation, then you only need gradients!

Gradient approximation for VI

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

$$\begin{aligned}\mathcal{L}_{\text{VI}}(q_\phi) &= \mathbb{E}_q [\log p(\mathbf{x}, \mathbf{z})] + \mathbb{H}[q_\phi(\mathbf{z}|\mathbf{x})] \\ &= \mathbb{E}_\pi [\log p(\mathbf{x}, \mathbf{f}_\phi(\epsilon, \mathbf{x}))] + \mathbb{H}[q_\phi(\mathbf{z}|\mathbf{x})] \quad // \text{ reparam. trick}\end{aligned}$$

If you use gradient descent for optimisation, then you only need gradients!

The gradient of the variational lower-bound:

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

Gradient approximation for VI

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

$$\begin{aligned}\mathcal{L}_{\text{VI}}(q_\phi) &= \mathbb{E}_q [\log p(\mathbf{x}, \mathbf{z})] + \mathbb{H}[q_\phi(\mathbf{z}|\mathbf{x})] \\ &= \mathbb{E}_\pi [\log p(\mathbf{x}, \mathbf{f}_\phi(\epsilon, \mathbf{x}))] + \mathbb{H}[q_\phi(\mathbf{z}|\mathbf{x})] \quad // \text{ reparam. trick}\end{aligned}$$

If you use gradient descent for optimisation, then you only need gradients!

The gradient of the variational lower-bound:

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

The gradient of the entropy term:

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

this term is 0

It remains to approximate $\nabla_{\mathbf{z}} \log q(\mathbf{z}|\mathbf{x})$!

Stein gradient estimator

Goal: approximate $\nabla_{\mathbf{x}} \log q(\mathbf{x})$ for a given distribution $q(\mathbf{x})$

Stein gradient estimator

Goal: approximate $\nabla_{\mathbf{x}} \log q(\mathbf{x})$ for a given distribution $q(\mathbf{x})$

Stein's identity:

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**:

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

Stein gradient estimator

Goal: approximate $\nabla_x \log q(\mathbf{x})$ for a given distribution $q(\mathbf{x})$

Stein's identity:

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)

Goal: approximate $\nabla_{\mathbf{x}} \log q(\mathbf{x})$ for a given distribution $q(\mathbf{x})$

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. Monte Carlo (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)

Goal: approximate $\nabla_{\mathbf{x}} \log q(\mathbf{x})$ for a given distribution $q(\mathbf{x})$

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

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

2. Rewrite the MC equations in matrix forms: denoting

$$\begin{aligned} \mathbf{H} &= (\mathbf{h}(\mathbf{x}^1), \dots, \mathbf{h}(\mathbf{x}^K)), \quad \overline{\nabla_{\mathbf{x}} \mathbf{h}} = \frac{1}{K} \sum_{k=1}^K \nabla_{\mathbf{x}^k} \mathbf{h}(\mathbf{x}^k), \\ \mathbf{G} &:= (\nabla_{\mathbf{x}^1} \log q(\mathbf{x}^1), \dots, \nabla_{\mathbf{x}^K} \log q(\mathbf{x}^K))^T, \end{aligned}$$

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

Stein gradient estimator (kernel based)

Goal: approximate $\nabla_{\mathbf{x}} \log q(\mathbf{x})$ for a given distribution $q(\mathbf{x})$

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

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)

Goal: approximate $\nabla_{\mathbf{x}} \log q(\mathbf{x})$ for a given distribution $q(\mathbf{x})$

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

3. Now solve a ridge regression problem:

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

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

with

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

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

Stein gradient estimator (kernel based)

Kernelized Stein Discrepancy:

$$\mathcal{S}^2(q, \hat{q}) = \mathbb{E}_{\mathbf{x}, \mathbf{x}' \sim q} [\hat{\mathbf{g}}(\mathbf{x})^\top \mathcal{K}_{\mathbf{xx}'} \hat{\mathbf{g}}(\mathbf{x}') + \hat{\mathbf{g}}(\mathbf{x})^\top \nabla_{\mathbf{x}'} \mathcal{K}_{\mathbf{xx}'} + \nabla_{\mathbf{x}} \mathcal{K}_{\mathbf{xx}'}^\top \hat{\mathbf{g}}(\mathbf{x}') + \text{Tr}(\nabla_{\mathbf{x}, \mathbf{x}'} \mathcal{K}_{\mathbf{xx}'})],$$

$$\mathbf{g}(\mathbf{x}) = \nabla_{\mathbf{x}} \log q(\mathbf{x}), \quad \hat{\mathbf{g}}(\mathbf{x}) = \nabla_{\mathbf{x}} \log \hat{q}(\mathbf{x}), \quad \mathcal{K}_{\mathbf{xx}'} = \mathcal{K}(\mathbf{x}, \mathbf{x}').$$

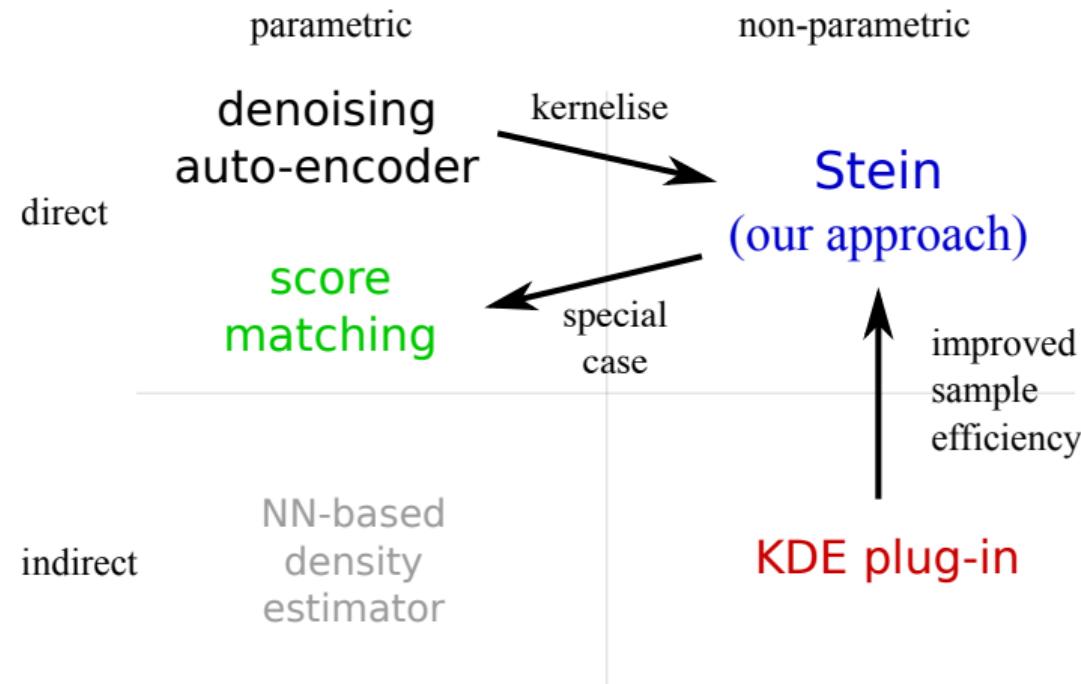
One can show that the V-statistic of KSD is

$$\mathcal{S}_V^2(q, \hat{q}) = \frac{1}{K^2} \text{Tr}(\hat{\mathbf{G}}^\top \mathbf{K} \hat{\mathbf{G}} + 2\hat{\mathbf{G}}^\top \langle \nabla, \mathbf{K} \rangle) + C$$

This means

$$\hat{\mathbf{G}}_V^{Stein} = \arg \min_{\hat{\mathbf{G}} \in \mathbb{R}^{K \times d}} \mathcal{S}_V^2(q, \hat{q}) + \frac{\eta}{K^2} \|\hat{\mathbf{G}}\|_F^2$$

Comparisons to existing approaches



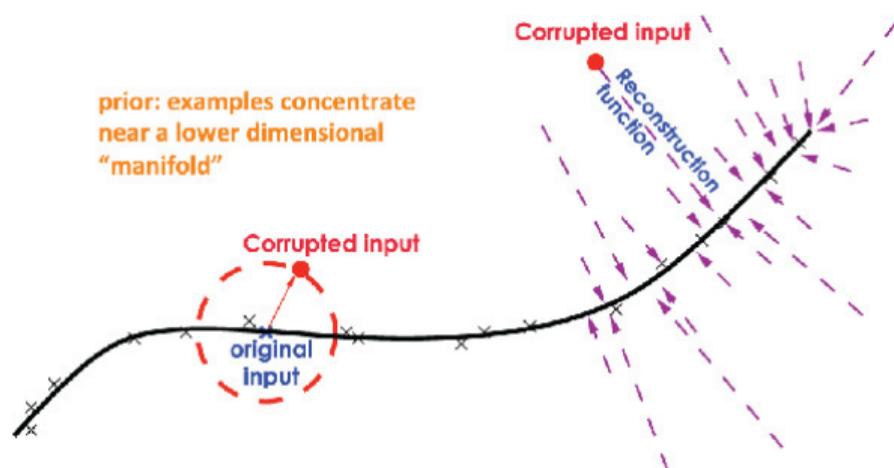
KDE plug-in estimator: Singh (1977)

Score matching estimator: Hyvärinen (2005), Sasaki et al. (2014), Strathmann et al. (2015)

Denoising auto-encoder: Vincent et al. (2008), Alain and Bengio (2014)

Comparisons to existing approaches

Compare to denoising auto-encoder (DAE):



- DAE: for $x \sim q(x)$, denoise $\hat{x} = x + \sigma\epsilon$ to x
(by min. ℓ_2 loss, $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$)
- When $\sigma \rightarrow 0$,
 $\text{DAE}^*(\hat{x}) \approx x + \sigma^2 \nabla_x \log q(x)$
 - unstable estimate: depends on σ
 - + functional gradient in RKHS:
 $\|\nabla \text{DAE loss}\|_{\mathcal{H}}^2 \propto \text{KSD}$

Vincent et al. (2008), Alain and Bengio (2014)
with Wenbo Gong and José Miguel Hernández-Lobato

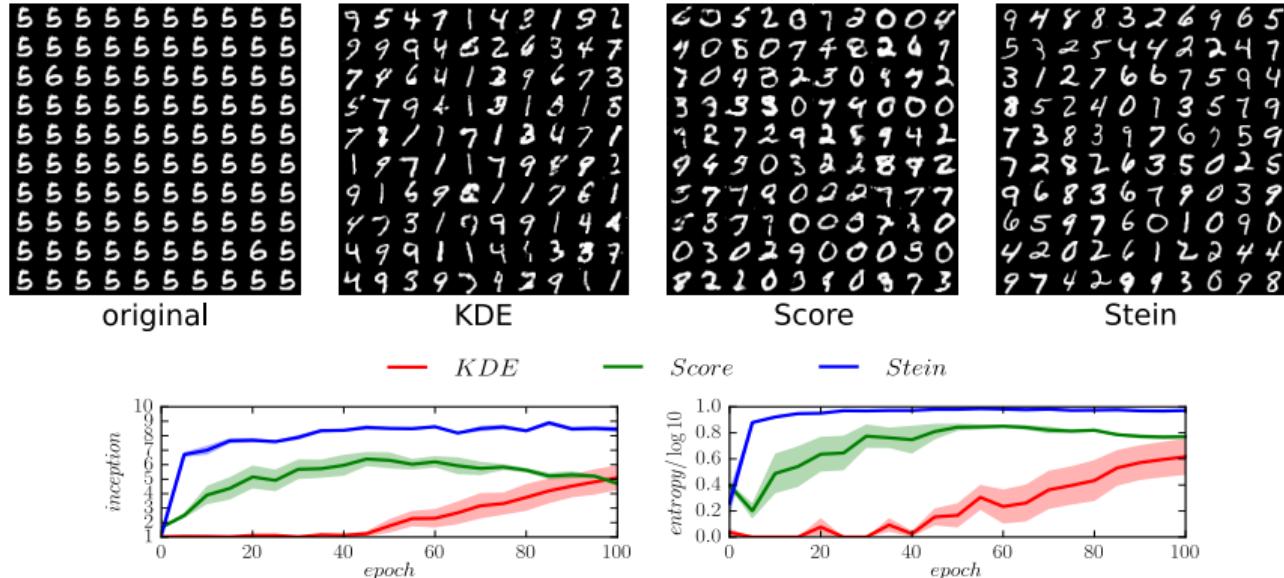
Example: entropy regularised GANs

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

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

- $\mathcal{L}_{\text{gen}}(p_{\text{gen}})$ is the generator loss of your favourite GAN method
- Again the gradient of $\mathbb{H}[p_{\text{gen}}]$ is approximated by the gradient estimators

Example: entropy regularised GANs



Significant improvement on sample diversity with the Stein approach

BEGAN: Berthelot et al. (2017)

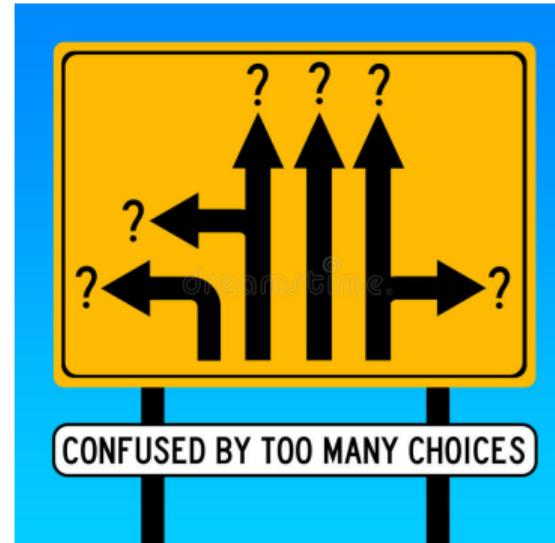
Meta learning for posterior samplers

Many existing posterior samplers in the literature...

- Which sampler should I use?
- How do I tune the hyper-parameters?

Learn a sampler from data!

- Want a general solution for **similar tasks**
- Train on low-dim, generalise to high-dim



Salimans et al. (2015), Song et al. (2017), Levy et al. (2018)

Learning to learn

Meta-learning for SG-MCMC

- Define a sampler with parameters ϕ :

$$\mathbf{z}_{t+1} = \mathbf{z}_t - \eta \mathbf{f}_\phi(\mathbf{z}_t, H(\cdot), \boldsymbol{\epsilon}), \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

- Run it on some training distributions $\pi(\mathbf{z}) \propto \exp[-H(\mathbf{z})]$, provide learning signals to train ϕ
- Once learned, apply this sampler to test distributions

Andrychowicz et al. (2016), Li and Malik (2017), Wichrowska et al. (2017), Li and Turner (2018)

The complete framework: Ma et al. NIPS 2015

- Itô diffusion

$$d\mathbf{z} = \mu(\mathbf{z})dt + \sqrt{2\mathbf{D}(\mathbf{z})}dW(t) \quad (1)$$

- To make sure $\pi(\mathbf{z}) \propto \exp[-H(\mathbf{z})]$ is a stationary distribution:

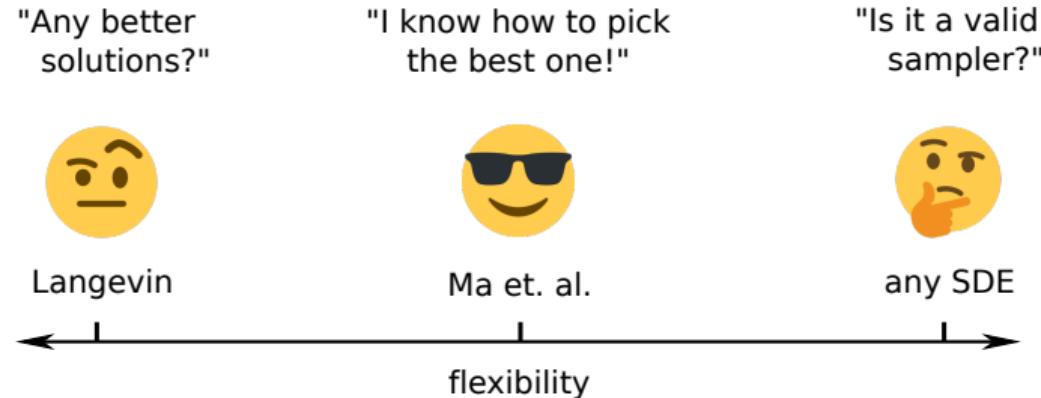
$$\mu(\mathbf{z}) = -[\mathbf{D}(\mathbf{z}) + \mathbf{Q}(\mathbf{z})]\nabla_{\mathbf{z}}H(\mathbf{z}) + \boldsymbol{\Gamma}(\mathbf{z}), \quad \boldsymbol{\Gamma}(\mathbf{z})_i = \sum_{j=1}^d \frac{\partial}{\partial z_j} [\mathbf{D}_{ij}(\mathbf{z}) + \mathbf{Q}_{ij}(\mathbf{z})] \quad (2)$$

- $\mathbf{D}(\mathbf{z})$: diffusion matrix, PSD
- $\mathbf{Q}(\mathbf{z})$: curl matrix, skew-symmetric
- $\boldsymbol{\Gamma}(\mathbf{z})$: correction vector

Ma et al. (2015) completeness result: under some mild conditions

“Any Itô diffusion that has the unique stationary $\pi(\mathbf{z})$ is governed by (1)+(2)”

The complete framework: Ma et al. NIPS 2015



- Searching the best sampler within the **complete** framework:
 - Guaranteed to be correct
 - Retains the most flexibility
 - Only needs to learn **how to parameterise $D(z)$ and $Q(z)$** matrices!

Our recipe: dynamics design

- **Goal:** train an SG-MCMC sampler to sample from $p(\theta|\mathcal{D}) \propto \exp[-U(\theta)]$
- We augment the state space with **momentum** variable \mathbf{p} :

$$\mathbf{z} = (\theta, \mathbf{p}), \quad \pi(\mathbf{z}) \propto \exp[-H(\mathbf{z})], \quad H(\mathbf{z}) = U(\theta) + \frac{1}{2}\mathbf{p}^\top \mathbf{p}$$

- Recall the complete recipe

$$d\mathbf{z} = -[\mathbf{D}(\mathbf{z}) + \mathbf{Q}(\mathbf{z})]\nabla_{\mathbf{z}} H(\mathbf{z})dt + \boldsymbol{\Gamma}(\mathbf{z})dt + \sqrt{2}dW(t)$$

Our recipe: dynamics design

- Our recipe:

$$\mathbf{Q}(\mathbf{z}) = \begin{bmatrix} \mathbf{0} & -\mathbf{Q}_f(\mathbf{z}) \\ \mathbf{Q}_f(\mathbf{z}) & \mathbf{0} \end{bmatrix}, \quad \mathbf{D}(\mathbf{z}) = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{D}_f(\mathbf{z}) \end{bmatrix}, \quad \boldsymbol{\Gamma}(\mathbf{z}) = \begin{bmatrix} \boldsymbol{\Gamma}_{\theta}(\mathbf{z}) \\ \boldsymbol{\Gamma}_{\rho}(\mathbf{z}) \end{bmatrix}$$

$$\mathbf{Q}_f(\mathbf{z}) = \text{diag}[\mathbf{f}_{\phi_Q}(\mathbf{z})], \quad \mathbf{D}_f(\mathbf{z}) = \text{diag}[\alpha \mathbf{f}_{\phi_Q}(\mathbf{z}) \odot \mathbf{f}_{\phi_Q}(\mathbf{z}) + \mathbf{f}_{\phi_D}(\mathbf{z}) + c], \quad \alpha, c > 0$$

- Resulting update rules (rearrange terms & discretise & stochastic gradient):

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \overbrace{\eta \mathbf{Q}_f(\mathbf{z}_t) \mathbf{p}_t}^{\text{momentum SGD}} + \overbrace{\eta \boldsymbol{\Gamma}_{\theta}(\mathbf{z}_t)}^{\text{correction}}$$

$$\mathbf{p}_{t+1} = \mathbf{p}_t - \underbrace{\eta \mathbf{D}_f(\mathbf{z}_t) \mathbf{p}_t}_{\text{friction}} - \eta \mathbf{Q}_f(\mathbf{z}_t) \nabla_{\boldsymbol{\theta}_t} \tilde{U}(\boldsymbol{\theta}_t) + \eta \boldsymbol{\Gamma}_{\rho}(\mathbf{z}_t) + \sqrt{\Sigma(\mathbf{z}_t)} \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$\Sigma(\mathbf{z}_t) = 2\eta \mathbf{D}_f(\mathbf{z}_t) - \eta^2 \mathbf{Q}_f(\mathbf{z}_t) \mathbf{B}(\boldsymbol{\theta}_t) \mathbf{Q}_f(\mathbf{z}_t), \quad \mathbf{B}(\boldsymbol{\theta}_t) = \mathbb{V}[\nabla_{\boldsymbol{\theta}_t} \tilde{U}(\boldsymbol{\theta}_t)]$$

Our recipe: dynamics design

Designing $\mathbf{f}_{\phi_Q}(\mathbf{z})$ (responsible for the drift):
the i^{th} element is defined as

$$\mathbf{f}_{\phi_Q,i}(\mathbf{z}) = \beta + f_{\phi_Q}(\tilde{U}(\boldsymbol{\theta}), p_i)$$

- We want $\mathbf{f}_{\phi_Q}(\mathbf{z})$ to depend on the energy landscape:
 - Fast traversal through low-density regions
 - Better exploration in high-density regions
- But we don't want $\Gamma_\theta(\mathbf{z})$ to be too expensive!
(using $\nabla_\theta U(\theta)$ as input here leads to an extra term $\langle \nabla, \nabla_\theta U(\theta) \rangle$ in $\Gamma_\theta(\mathbf{z})$)

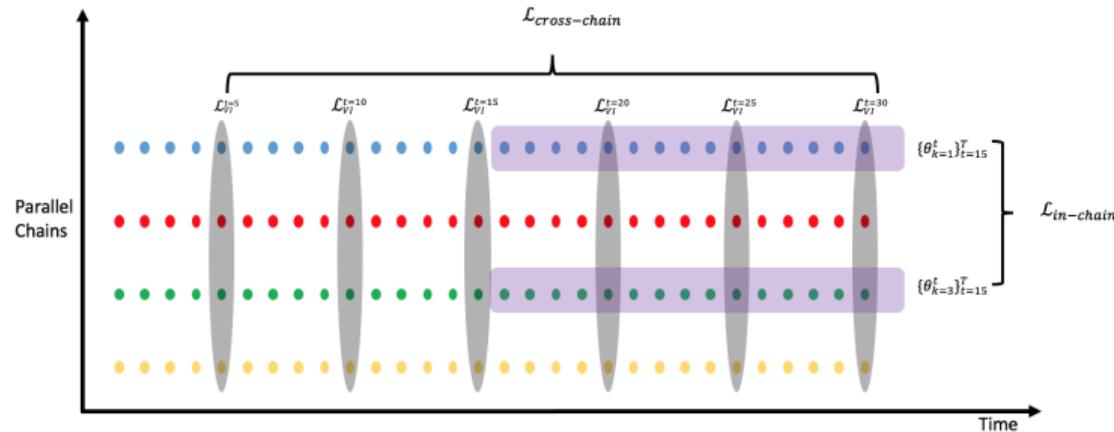
Our recipe: dynamics design

Designing $\mathbf{f}_{\phi_D}(\mathbf{z})$ (responsible for friction):
the i^{th} element is defined as

$$\mathbf{f}_{\phi_D,i}(\mathbf{z}) = f_{\phi_D}(\tilde{U}(\boldsymbol{\theta}), p_i, \partial_{\theta_i} \tilde{U}(\boldsymbol{\theta}))$$

- $\Gamma_p(z)$ only requires computing $\nabla_p D_f(z)$
- ...so we can use the gradient information $\nabla_{\theta} U(\boldsymbol{\theta})$
- prevent overshoot by “comparing” p and $\nabla_{\theta} U(\boldsymbol{\theta})$

Our recipe: loss function design



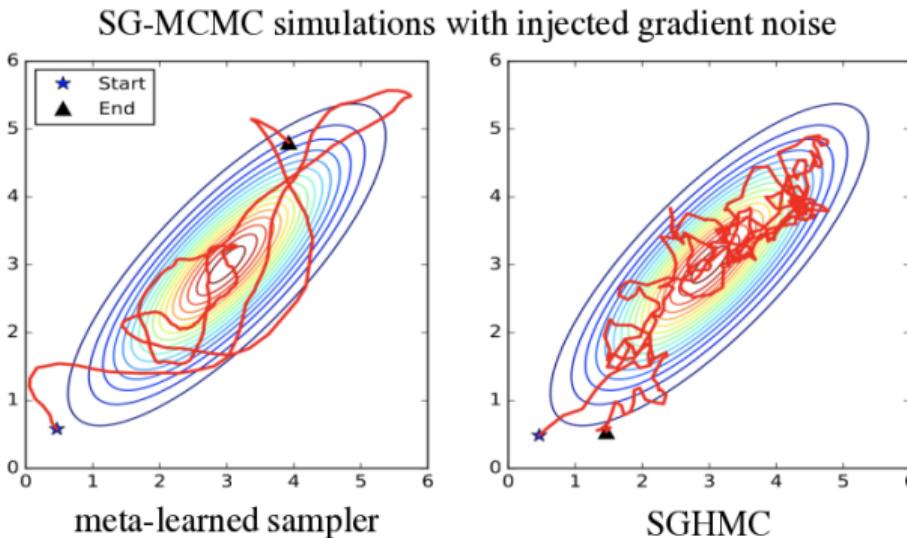
Use KL divergence $\text{KL}[q(\theta)||p(\theta|\mathcal{D})]$ to define loss.

Define $q(\theta)$ implicitly: run parallel chains for several steps, then

- Cross-chain loss: at time t , collect samples across chains
- In-chain loss: for each chain, collect samples by *thinning*

Gradient of KL approximated by Stein gradient estimator!

A toy example



- trained on **factorised** Gaussians, tested on **correlated** Gaussians
- manually injected Gaussian noise to the gradients
(and assume we don't know noise variance $B(\theta)$)

Bayesian NN on MNIST

Goal: sample from BNN posterior

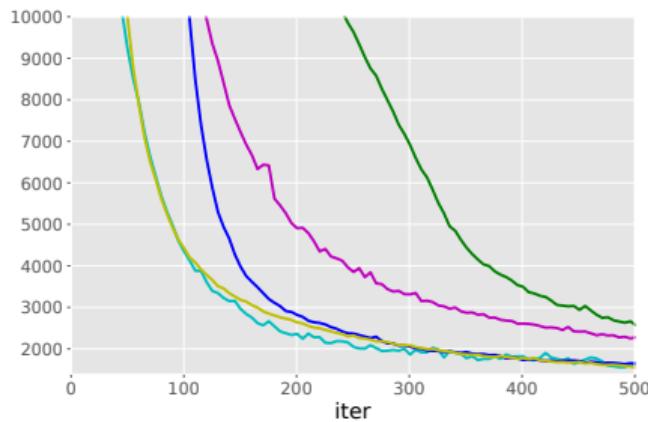
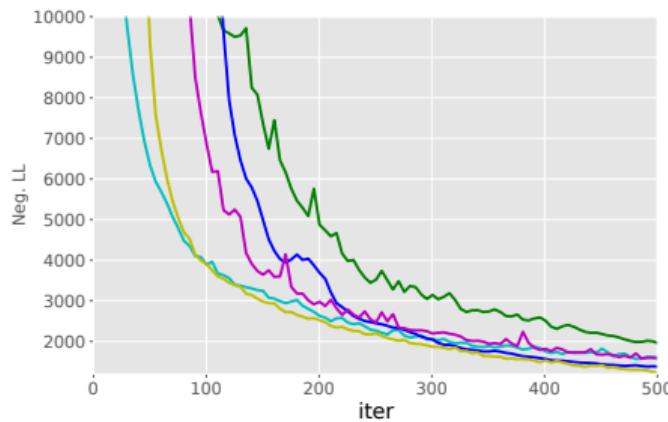
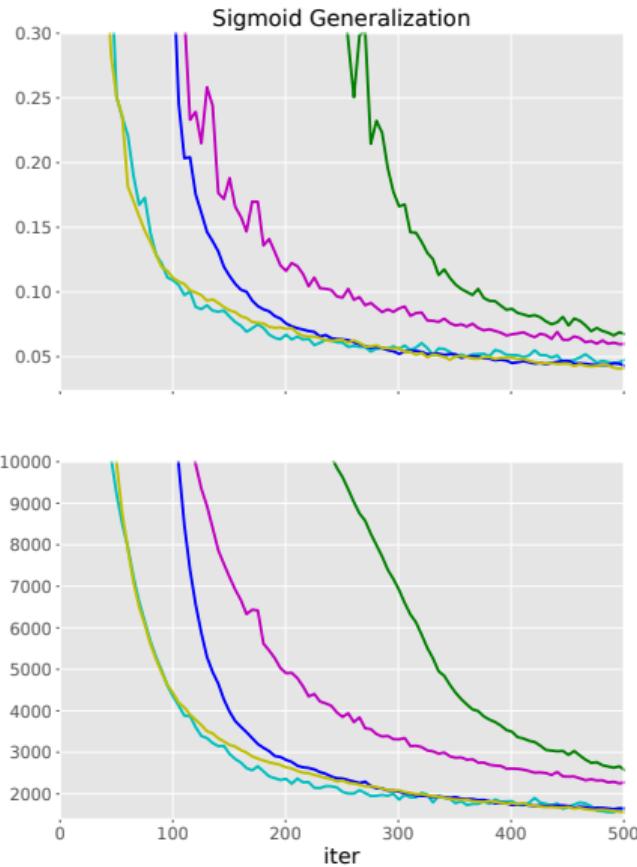
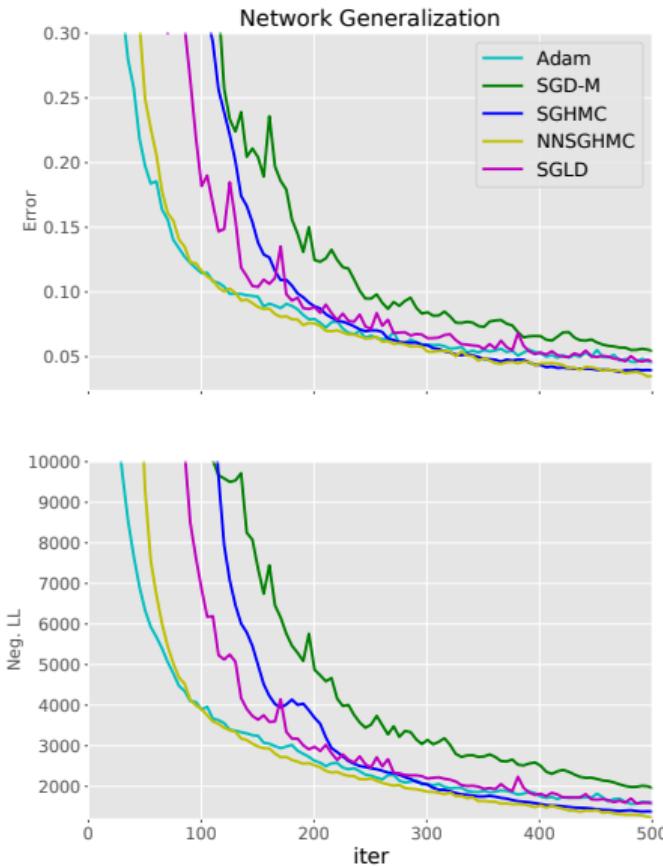
Training: meta sampler trained to sample from the posterior of a BNN
(1-hidden layer, 20 hidden units, ReLU)

Three generalisation tests:

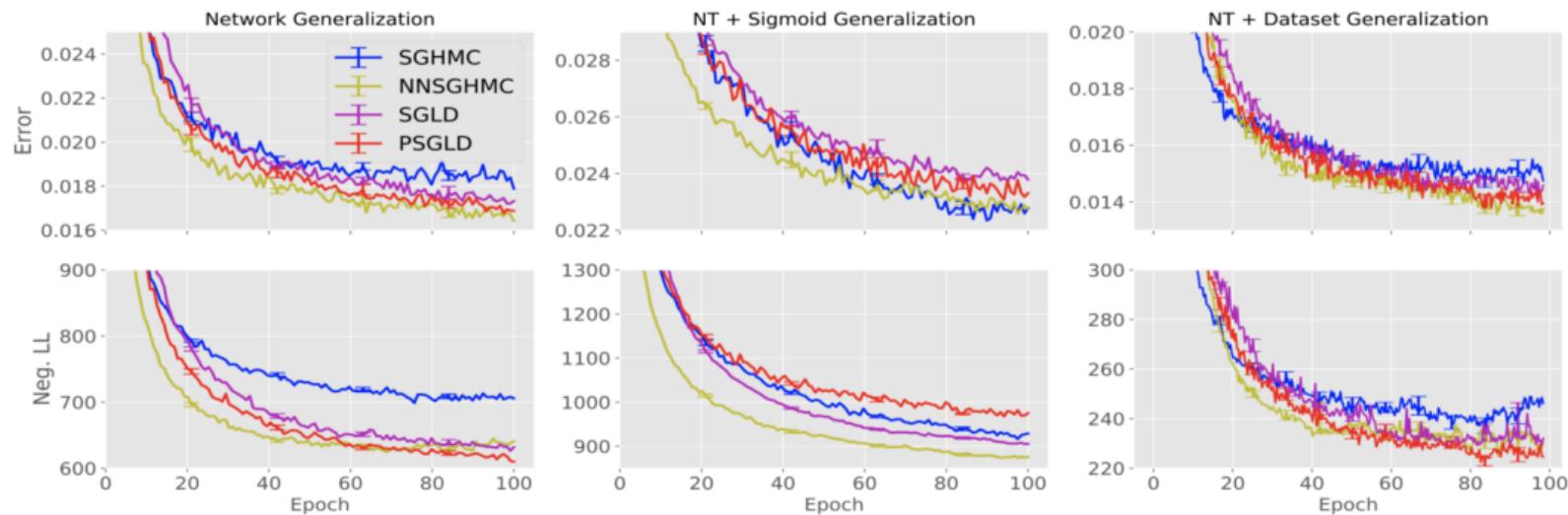
- to **bigger** network architecture: 2-hidden layer MLP (**40** units, ReLU)
- to **different** activation function: 1-hidden layer MLP (20 units, **Sigmoid**)
- to **different** dataset: train on **MNIST 0-4**, test on **MNIST 5-9**

Also consider **long-time** horizon generalisation

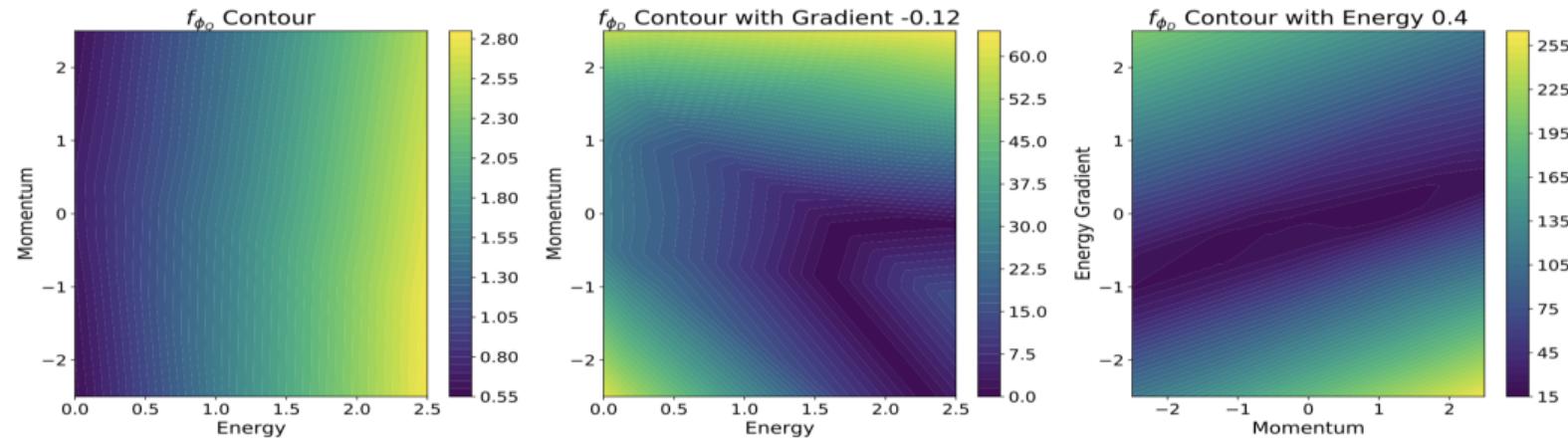
Bayesian NN on MNIST: speed improvements



Bayesian NN on MNIST: long-time generalisation



Bayesian NN on MNIST: understanding the learned sampler

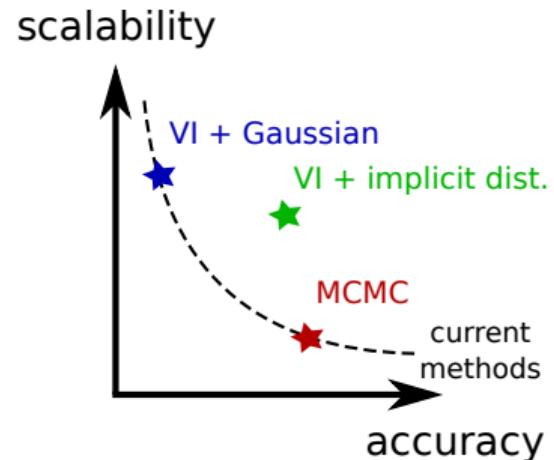


$$\mathbf{Q}_f(\mathbf{z}) = \text{diag}[\mathbf{f}_{\phi_Q}(\mathbf{z})], \quad \mathbf{D}_f(\mathbf{z}) = \text{diag}[\alpha \mathbf{f}_{\phi_Q}(\mathbf{z}) \odot \mathbf{f}_{\phi_Q}(\mathbf{z}) + \mathbf{f}_{\phi_D}(\mathbf{z}) + c]$$

- \mathbf{f}_{ϕ_Q} (left): nearly linear wrt. energy (fast traversal, better exploration)
- \mathbf{f}_{ϕ_D} (middle): decrease friction around high energy regions
- \mathbf{f}_{ϕ_D} (right): increase friction when gradient & momentum “disagree” (prevent overshoot)

Stein gradient estimator: summary

- We derived a non-parametric gradient estimator
- We used the gradient estimator for entropy-regularised GANs
- We explored meta-learning for approximate inference



Related work by colleagues

Doucet et al. (2013). Derivative-Free Estimation of the Score Vector and Observed Information Matrix with Application to State-Space Models. ArXiv 1304.5768

Shi et al. (2018). A Spectral Approach to Gradient Estimation for Implicit Distributions. ICML 2018

Song et al. (2019). Sliced Score Matching: A Scalable Approach to Density and Score Estimation. UAI 2019

Andrew Duncan's talk. Minimum Stein discrepancy estimators. [This workshop](#)
Spectral Estimators for Gradient Fields of Log-Densities. [This workshop](#)



Thank you!

Y. Li and R.E. Turner. Gradient Estimators for Implicit Models. ICLR 2018

W. Gong*, Y. Li* and J.M. Hernández-Lobato. Meta learning for stochastic gradient MCMC. ICLR 2019.