

# Meta-learning for stochastic gradient MCMC

---

Yingzhen Li

University of Cambridge → Microsoft Research Cambridge

Joint work with Wenbo Gong & José Miguel Hernández-Lobato (U Cambridge)  
arXiv 1806.04522

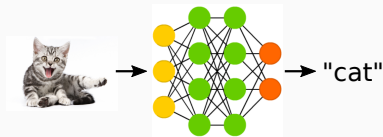


# Bayesian neural networks 101

**Goal:** classifying different types of cats from images

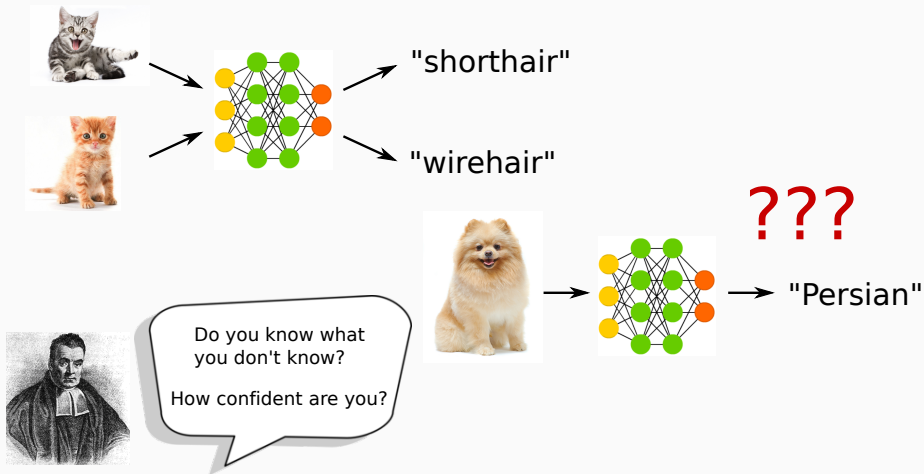
- $\mathbf{x}$ : input images;  $\mathbf{y}$ : output label
- build a neural network (with param.  $\theta$ ):  
 $\hat{\mathbf{y}} = \text{NN}_{\theta}(\mathbf{x})$
- find the best parameter  $\theta$  given a dataset  $\mathcal{D} = \{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^N$ :

$$\theta^* = \arg \max_{\theta} \sum_{n=1}^N \log p(\mathbf{y}_n | \mathbf{x}_n, \theta) + \log p(\theta)$$



Maximum a posteriori (MAP)

# Bayesian neural networks 101



**Bayesian inference:** given some function  $F(\theta)$ , want  $\mathbb{E}_{p(\theta|\mathcal{D})} [F(\theta)]$



$$P(\theta|\mathcal{D}) = \frac{P(\mathcal{D}|\theta)P(\theta)}{P(\mathcal{D})}$$

- predictive mean

$$\hat{\mathbf{y}}_{\text{mean}} = \mathbb{E}_{p(\theta|\mathcal{D})} [\text{NN}_{\theta}(\mathbf{x})]$$

- predictive distribution

$$p(\mathbf{y}|\mathbf{x}, \mathcal{D}) = \mathbb{E}_{p(\theta|\mathcal{D})} [p(\mathbf{y}|\mathbf{x}, \theta)]$$

- evaluate posterior

$$p(\theta \in A|\mathcal{D}) = \mathbb{E}_{p(\theta|\mathcal{D})} [\delta_A]$$

**Bayesian inference:** given some function  $F(\theta)$ , want  $\mathbb{E}_{p(\theta|\mathcal{D})} [F(\theta)]$



$$P(\theta|\mathcal{D}) = \frac{P(\mathcal{D}|\theta)P(\theta)}{P(\mathcal{D})}$$

- Monte Carlo estimation:

$$\mathbb{E}_{p(\theta|\mathcal{D})} [F(\theta)] \approx \frac{1}{K} \sum_{k=1}^K F(\theta^k)$$

$$\theta^k \sim p(\theta|\mathcal{D}) \quad (\text{intractable})$$

- Stochastic gradient MCMC (SG-MCMC):  
efficient ways to (approximately) draw  
samples from  $p(\theta|\mathcal{D})$

- The MAP problem can be rewritten as

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} U(\boldsymbol{\theta}),$$

$$p(\boldsymbol{\theta}|\mathcal{D}) \propto \exp[-U(\boldsymbol{\theta})], \quad -U(\boldsymbol{\theta}) = \sum_{n=1}^N \log p(\mathbf{y}_n|\mathbf{x}_n, \boldsymbol{\theta}) + \log p(\boldsymbol{\theta})$$

## From SGD to SG-MCMC

- The MAP problem can be rewritten as

$$\theta^* = \arg \min_{\theta} U(\theta),$$

$$p(\theta|\mathcal{D}) \propto \exp[-U(\theta)], \quad -U(\theta) = \sum_{n=1}^N \log p(\mathbf{y}_n|\mathbf{x}_n, \theta) + \log p(\theta)$$

- In the MAP problem, we find  $\theta^*$  by gradient descent

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} U(\theta_t),$$

$$-\nabla_{\theta} U(\theta) = \underbrace{\sum_{n=1}^N \nabla_{\theta} \log p(\mathbf{y}_n|\mathbf{x}_n, \theta) + \nabla_{\theta} \log p(\theta)}_{\text{full gradient of LL}}, \quad \{(\mathbf{x}_m, \mathbf{y}_m)\}_{m=1}^M \sim \mathcal{D}^M$$

# From SGD to SG-MCMC

- The MAP problem can be rewritten as

$$\theta^* = \arg \min_{\theta} U(\theta),$$

$$p(\theta|\mathcal{D}) \propto \exp[-U(\theta)], \quad -U(\theta) = \sum_{n=1}^N \log p(\mathbf{y}_n|\mathbf{x}_n, \theta) + \log p(\theta)$$

- In the MAP problem, we find  $\theta^*$  by stochastic gradient descent (for big data)

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} \tilde{U}(\theta_t),$$

$$-\nabla_{\theta} \tilde{U}(\theta) = \underbrace{\frac{N}{M} \sum_{m=1}^M \nabla_{\theta} \log p(\mathbf{y}_m|\mathbf{x}_m, \theta)}_{\text{stochastic gradient of LL}} + \nabla_{\theta} \log p(\theta), \quad \{(\mathbf{x}_m, \mathbf{y}_m)\}_{m=1}^M \sim \mathcal{D}^M$$



# From SGD to SG-MCMC

- In the Bayesian inference problem, we need to (approximately) draw  $\theta \sim p(\theta|\mathcal{D})$
- Big data: Stochastic gradient Langevin dynamics (SGLD)

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} \tilde{U}(\theta_t) + \sqrt{2\eta} \epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Welling and Teh (2011), Chen et al. (2014), Li et al. (2016), Chen et al. (2016)

# From SGD to SG-MCMC

- In the Bayesian inference problem, we need to (approximately) draw  $\theta \sim p(\theta|\mathcal{D})$
- Big data: Stochastic gradient Langevin dynamics (SGLD)

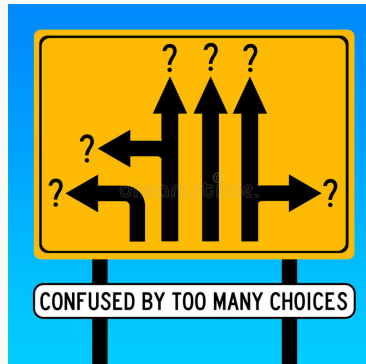
$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} \tilde{U}(\theta_t) + \sqrt{2\eta} \epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

- SGLD = SGD + properly scaled Gaussian noise
- Other optimisation algorithms can be transformed into SG-MCMC samplers:
  - SGD + momentum  $\rightarrow$  SGHMC; RMSprop  $\rightarrow$  preconditioned SGLD; Adam  $\rightarrow$  Santa

Welling and Teh (2011), Chen et al. (2014), Li et al. (2016), Chen et al. (2016)

# “I’m bored of tuning my optimiser & sampler”

- Which **SG-MCMC** algorithm should I use?
- How do I tune the hyper-parameters?



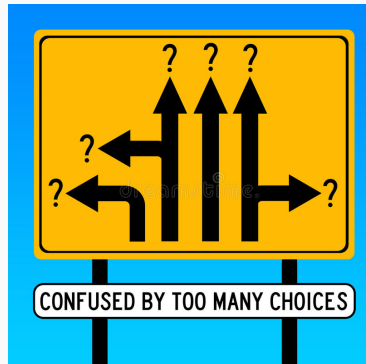
Salimans et al. (2015), Song et al. (2017), Levy et al. (2018)

# “I’m bored of tuning my optimiser & sampler”

- Which **SG-MCMC** algorithm should I use?
- How do I tune the hyper-parameters?

Learn it 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)

Meta-learning for **optimisers**:

- Define **an optimiser** with parameters  $\phi$ :

$$\mathbf{z}_{t+1} = \mathbf{z}_t - \mathbf{f}_{\phi}(\mathbf{z}_t, H(\cdot))$$

- Run it on some training **objective functions**  $H(\mathbf{z})$ ,  
provide learning signals to train  $\phi$
- Once learned, apply this **optimiser** to test **objective functions**

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

Meta-learning for SG-MCMC: can we just naively

- Define a sampler with parameters  $\phi$ :

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

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

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

Meta-learning for SG-MCMC: can we just naively

- Define a **sampler** with parameters  $\phi$ :

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

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

**Not quite yet!** We need to make sure it is a valid sampler!

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

To sample from  $\pi(\mathbf{z}) \propto \exp[-H(\mathbf{z})]$ :

- Let's take the step-size  $\eta \rightarrow 0$  and use exact gradient:

$$d\mathbf{z} = -\nabla_{\mathbf{z}} H(\mathbf{z})dt + \sqrt{2}dW(t) \quad (\text{Langevin dynamics})$$

- $W(t)$  is a Wiener process  
(think about  $dW(t)$  as some Gaussian noise with variance  $dt$ )
- Langevin dynamics is a special case of **Itô diffusion**

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



# The complete framework: Ma et al. NIPS 2015

- Itô diffusion

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

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

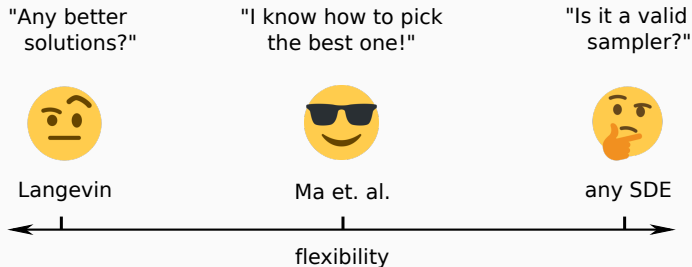
$$\mu(z) = -[D(z) + Q(z)]\nabla_z H(z) + \Gamma(z), \quad \Gamma(z)_i = \sum_{j=1}^d \frac{\partial}{\partial z_j} [D_{ij}(z) + Q_{ij}(z)] \quad (2)$$

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

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

“Any Itô diffusion that has the unique stationary  $\pi(z)$  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  $p(\boldsymbol{\theta}|\mathcal{D}) \propto \exp[-U(\boldsymbol{\theta})]$
- We augment the state space with **momentum** variable  $\mathbf{p}$ :

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

- 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:

$$\mathbf{Q}(\mathbf{z}) = \begin{bmatrix} \mathbf{0} & -\mathbf{Q}_f(\mathbf{z}) \\ \mathbf{Q}_f(\mathbf{z}) & \mathbf{0} \end{bmatrix}, \quad \mathbf{D}(\mathbf{z}) = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{D}_f(\mathbf{z}) \end{bmatrix}, \quad \mathbf{\Gamma}(\mathbf{z}) = \begin{bmatrix} \mathbf{\Gamma}_\theta(\mathbf{z}) \\ \mathbf{\Gamma}_p(\mathbf{z}) \end{bmatrix}$$
$$\mathbf{Q}_f(\mathbf{z}) = \text{diag}[\mathbf{f}_{\phi_Q}(\mathbf{z})], \quad \mathbf{D}_f(\mathbf{z}) = \text{diag}[\alpha \mathbf{f}_{\phi_Q}(\mathbf{z}) \odot \mathbf{f}_{\phi_Q}(\mathbf{z}) + \mathbf{f}_{\phi_D}(\mathbf{z}) + c], \quad \alpha, c > 0$$

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

$$\begin{aligned} \theta_{t+1} &= \theta_t + \overbrace{\eta \mathbf{Q}_f(\mathbf{z}_t) \mathbf{p}_t}^{\text{momentum SGD}} + \overbrace{\eta \mathbf{\Gamma}_\theta(\mathbf{z}_t)}^{\text{correction}} \\ \mathbf{p}_{t+1} &= \mathbf{p}_t - \underbrace{\eta \mathbf{D}_f(\mathbf{z}_t) \mathbf{p}_t}_{\text{friction}} - \eta \mathbf{Q}_f(\mathbf{z}_t) \nabla_{\theta_t} \tilde{U}(\theta_t) + \eta \mathbf{\Gamma}_p(\mathbf{z}_t) + \sqrt{\Sigma(\mathbf{z}_t)} \epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ \Sigma(\mathbf{z}_t) &= 2\eta \mathbf{D}_f(\mathbf{z}_t) - \eta^2 \mathbf{Q}_f(\mathbf{z}_t) \mathbf{B}(\theta_t) \mathbf{Q}_f(\mathbf{z}_t), \quad \mathbf{B}(\theta_t) = \mathbb{V}[\nabla_{\theta_t} \tilde{U}(\theta_t)] \end{aligned}$$

# Our recipe: dynamics design

Designing  $\mathbf{f}_{\phi_Q}(\mathbf{z})$  (responsible for the drift):  
the  $i^{\text{th}}$  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  $\boldsymbol{\Gamma}_{\boldsymbol{\theta}}(\mathbf{z})$  to be too expensive!  
(using  $\nabla_{\boldsymbol{\theta}} U(\boldsymbol{\theta})$  as input here leads to an extra term  $\langle \nabla, \nabla_{\boldsymbol{\theta}} U(\boldsymbol{\theta}) \rangle$  in  $\boldsymbol{\Gamma}_{\boldsymbol{\theta}}(\mathbf{z})$ )

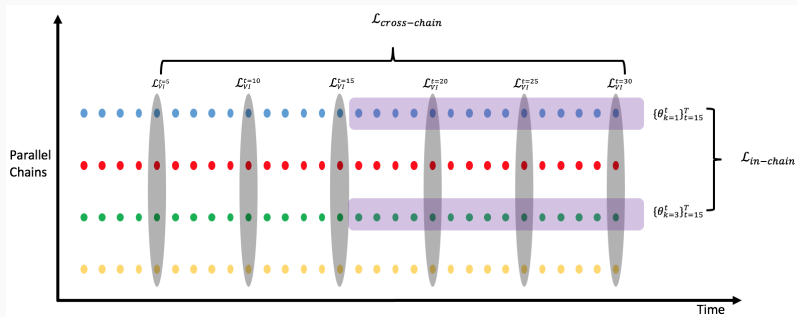
## Our recipe: dynamics design

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

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

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

# Our recipe: loss function design

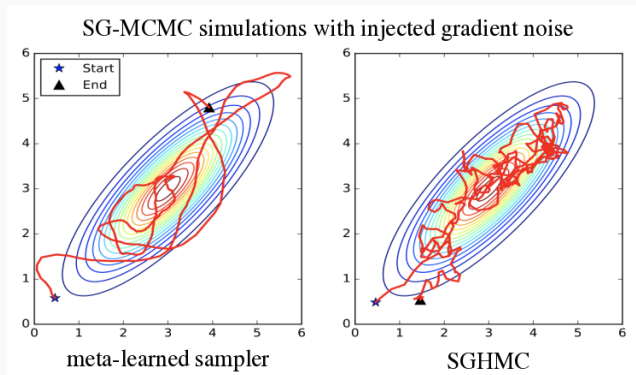


Use KL divergence  $\text{KL}[q(\theta)||p(\theta|\mathcal{D})]$  to define loss.

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*

## 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  $\mathbf{B}(\theta)$ )



**Goal:** sample from BNN posterior

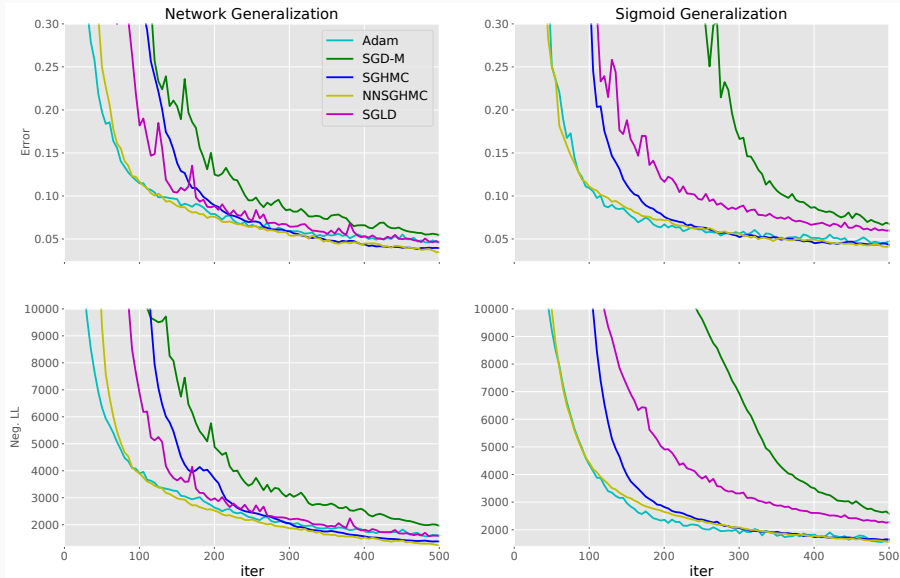
**Training:** meta sampler trained to sample from the posterior of a BNN (1-hidden layer, 20 hidden units, ReLU)

Three generalisation tests:

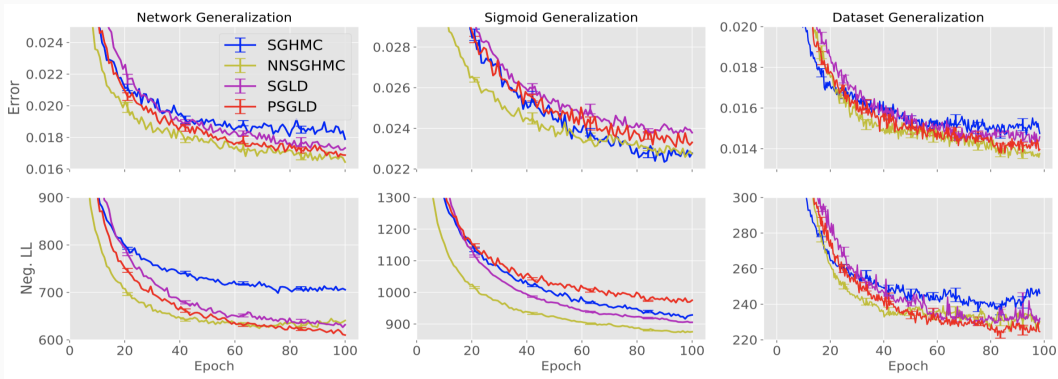
- to **bigger** network architecture: 2-hidden layer MLP (40 units, ReLU)
- to **different** activation function: 1-hidden layer MLP (20 units, Sigmoid)
- to **different** dataset: train on MNIST 0-4, test on MNIST 5-9

Also consider **long-time** horizon generalisation

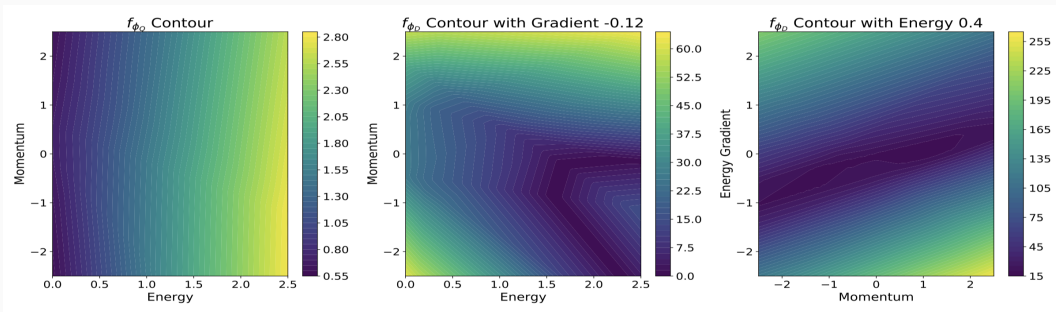
# Bayesian NN on MNIST: speed improvements



# Bayesian NN on MNIST: long-time generalisation



# Bayesian NN on MNIST: understanding the learned sampler



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

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

# Summary

MCMC and meta-learning can be friends:

- MCMC can be improved using meta-learning
- Meta-learning works better when searching in a theoretically sound framework



**Thank you!**

# Summary

MCMC and meta-learning can be friends:

- MCMC can be improved using meta-learning
- Meta-learning works better when searching in a theoretically sound framework

Future work:

- add in tempering, adaptive learning rate...
- meta-learn the Hamiltonian
- improving samplers for discrete variables



**Thank you!**

# References

- Welling and Teh (2011). Bayesian learning via stochastic gradient Langevin dynamics. ICML 2011
- Chen et al. (2014). Stochastic gradient Hamiltonian Monte Carlo. ICML 2014
- Li et al. (2016). Learning weight uncertainty with stochastic gradient MCMC for shape classification. CVPR 2016
- Chen et al. (2016) Bridging the gap between stochastic gradient MCMC and stochastic optimization. AISTATS 2016
- Salimans et al. (2015). Markov chain Monte Carlo and variational inference: Bridging the gap. ICML 2015
- Song et al. (2017). A-NICE-MC: Adversarial training for MCMC. NIPS 2017
- Levy et al. (2018). Generalizing Hamiltonian Monte Carlo with neural networks. ICLR 2018
- Andrychowicz et al. (2016). Learning to learn by gradient descent by gradient descent. NIPS 2016
- Li and Malik (2017). Learning to optimize. ICLR 2017
- Wichrowska et al. (2017). Learned optimizers that scale and generalize. ICML 2017
- Li and Turner (2018). Gradient estimators for implicit models. ICLR 2018
- Ma et al. (2015). A complete recipe for stochastic gradient MCMC. NIPS 2015