

# Inference and validation of post-Bayesian methods



Zheyang Shen  
Newcastle University

Lancaster University: CSML talk

# Papers

- ▶ *Prediction-centric uncertainty quantification via MMD.* Z. Shen, J. Knoblauch, S. Power, C. J. Oates.
- ▶ *A computable measure of suboptimality for entropy-regularised variational objectives.* C. Chazal, H. Kanagawa, Z. Shen, A. Korba, C. J. Oates.
- ▶ *Predictively oriented posteriors.* Y. McLatchie, B.-E. Cherief-Abdellatif, D. T. Frazier, J. Knoblauch.

# Optimization-centric probabilistic inference

Bayesian posterior and its generalization are minimizers of an *entropy-regularized objective*

$$P := \arg \min_{Q \in \mathcal{P}(\mathbb{R}^d)} \mathcal{J}(Q), \quad \mathcal{J}(Q) := \mathcal{L}(Q) + \text{KL}(Q \| Q_0).$$

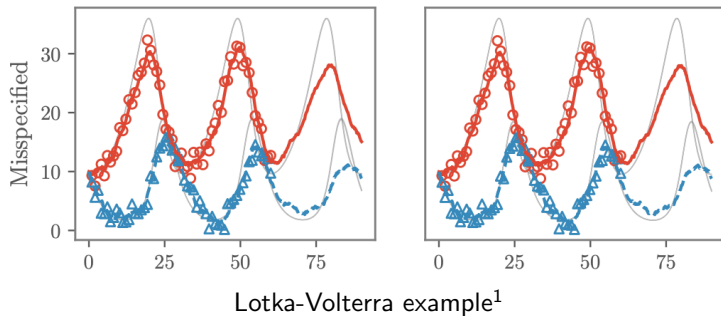
## Examples of $\mathcal{L}(Q)$

- ▶ Standard Bayes:  $-\sum_{i=1}^n \int \log p(y_i | \theta) \, dQ(\theta)$ .
- ▶ Generalized Bayes<sup>1</sup>:  $\sum_{i=1}^n \int \ell(y_i, \theta) \, dQ(\theta)$ .

<sup>1</sup> Bissiri et al. [2016]

# Why post-Bayes?

- ▶ Bayesian posterior in a misspecified model lack in parameter uncertainty;
- ▶ Generalized Bayes confronts aspects of model misspecification (e.g., outlier robustness), yet is still overconfident.



<sup>1</sup> Shen et al. [2025]

# Predictively-oriented (PrO) posteriors<sup>1</sup>

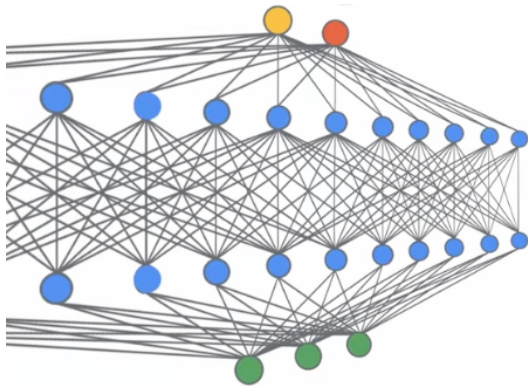
Generalized Bayes	Predictively-oriented
$\underbrace{\sum_{i=1}^n \int S(y_i, P_\theta) \, dQ(\theta)}_{\text{average fit}}$	$\underbrace{\sum_{i=1}^n S\left(y_i, \int P_\theta \, dQ(\theta)\right)}_{\text{predictive fit}}$

*What happens when the model is well-specified?*

McLatchie et al. [2025] illustrate that the PrO posterior still concentrates around the data-generating distribution (Theorem 1).

<sup>1</sup> McLatchie et al. [2025]

# Mean field neural networks (MFNNs)



The output of MFNN is “infinitely wide”

$$f_Q(x) = \int \Phi(x, \theta) \, dQ(\theta)$$

Training loss is a nonlinear  $\mathcal{L}(Q)$ :

$$\mathcal{L}(Q) = \sum_{i=1}^n \ell(y_i, f_Q(x_i)).$$

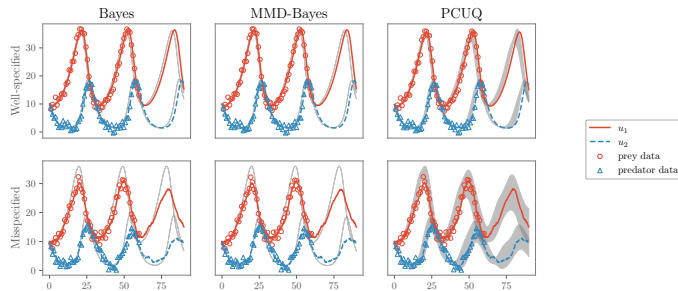
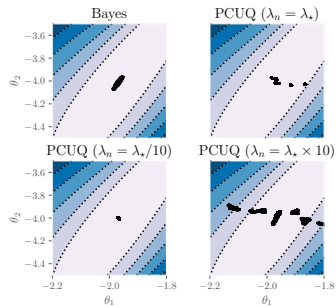
# Challenges of a nonlinear $\mathcal{L}(Q)$

$$P := \arg \min_{Q \in \mathcal{P}(\mathbb{R}^d)} \mathcal{L}(Q) + \text{KL}(Q \| Q_0).$$

- ▶  $P$  is identifiable up to normalization when  $\mathcal{L}(Q)$  is linear: applies to (generalized) Bayes;
- ▶ The nonlinearity of  $\mathcal{L}(Q)$  makes it so that  $P$  is only identifiable via optimization:  
mean-field Langevin dynamics evolve a set of interacting particles  $Q_N^t := \frac{1}{N} \sum_{i=1}^N \delta_{\theta_i^t}$ :

$$\theta_i^{t+1} = \theta_i^t + \epsilon [\nabla \log q_0(\theta_i^t) - \underbrace{\nabla_V \mathcal{L}(Q_N^t)(\theta_i^t)}_{\text{variational gradient}}] + \sqrt{2\epsilon} Z_i^t.$$

*Given a sample-based approximation  $Q_N := \frac{1}{N} \sum_{i=1}^N \delta_{\theta_i}$ , is it any good?*





## Revisiting Stein's identity

We would like to generalize Stein's discrepancy to validate how well a set of particles satisfies the minimizer of  $\mathcal{J}(Q)$ .

$$\forall \phi \in \mathcal{F}, \quad \mathbb{E}_{x \sim Q} [\langle \nabla \log p(x), \phi(x) \rangle + \langle \nabla, \phi(x) \rangle] = 0 \quad \Leftrightarrow \quad Q = P.$$

Noting that  $\mathcal{J}_{\text{Bayes}}(Q) = \text{KL}(Q \| P)$ , we have:

$$\begin{aligned} \mathbb{E}_{x \sim Q} [\langle \nabla \log p(x), \phi(x) \rangle + \langle \nabla, \phi(x) \rangle] &= \mathbb{E}_{x \sim Q} [\langle \nabla \log p(x) - \nabla \log q(x), \phi(x) \rangle] \\ &= \mathbb{E}_{x \sim Q} [\langle -\nabla_v \mathcal{J}_{\text{Bayes}}(Q)(x), \phi(x) \rangle] = - \left. \frac{d}{dt} \mathcal{J}_{\text{Bayes}}(Q_t^\phi) \right|_{t=0}. \end{aligned}$$

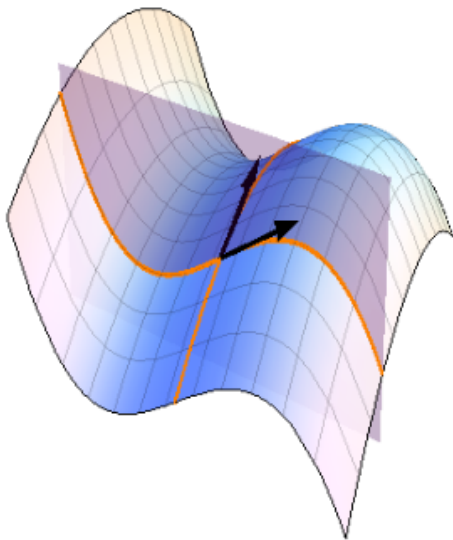
*Drawing any geodesic curve  $(Q_t^\phi)_{t \in (-\epsilon, \epsilon)}$  around  $Q$ , its time derivative w.r.t.  $\mathcal{J}_{\text{Bayes}}(Q_t)$  is always zero.*

# Measuring approximation quality

**Idea:** See “how well  $Q$  minimises  $\mathcal{J}$ ”

**Concretely:** If  $Q$  does **not** minimise  $\mathcal{J}$  then there is some “direction”  $\phi$  such that

$$\left. \frac{d}{dt} \mathcal{J}(Q_t^\phi) \right|_{t=0} < 0.$$



# Variational Gradients

**Variational gradients:** From the fundamental theorem of calculus

$$\left. \frac{d}{dt} \mathcal{J}(Q_t^\phi) \right|_{t=0} = \int \langle \nabla_V \mathcal{J}(Q)(x), \phi(x) \rangle dQ(x).$$

where the *variational gradient* is  $\nabla_V \mathcal{J}(Q)(x) := \nabla_x \mathcal{J}'(Q)(x)$  for each  $x \in \mathbb{R}^d$ .

(the *first variation*  $\mathcal{F}'(Q)$  is defined as  $\left. \frac{d}{d\epsilon} \mathcal{F}(Q + \epsilon\chi) \right|_{\epsilon=0} = \int \mathcal{F}'(Q) d\chi$ )

**Computing the variational gradient of  $\mathcal{J}$ :** Letting  $Q_0$  have a density  $q_0 > 0$ ,

$$\nabla_V \mathcal{J}(Q)(\theta) = \nabla_V \mathcal{L}(Q)(\theta) - (\nabla \log q_0)(\theta) + \underbrace{(\nabla \log q)(\theta)}_{\text{problematic}}$$

**Main idea:** Can still evaluate integrals of the variational gradient:

$$\int \underbrace{(\nabla \log q)(x)}_{\text{problematic}} \cdot \phi(x) dQ(x) = - \underbrace{\int (\nabla \cdot \phi)(x) dQ(x)}_{\text{fine}}.$$

# Gradient Discrepancy

**Gradient discrepancy:** For a given set  $\mathcal{F}$  of differentiable vector fields on  $\mathbb{R}^d$ , define the *gradient discrepancy* as

$$\text{GD}(Q) := \sup_{\substack{\phi \in \mathcal{F} \text{ s.t.} \\ (\mathcal{T}_Q \phi)_- \in \mathcal{L}^1(Q)}} \left| \int \mathcal{T}_Q \phi(x) \, dQ(x) \right|$$

where  $\mathcal{T}_Q \phi(x) := [(\nabla \log q_0)(x) - \nabla_V \mathcal{L}(Q)(x)] \cdot \phi(x) + (\nabla \cdot \phi)(x)$ .

**Example:** For  $\mathcal{L}'(Q)(x) = -\log p(y|x)$ , we recover (*Langevin*) *Stein discrepancy* [Gorham and Mackey, 2015]

$$\text{SD}(Q) := \sup_{\substack{\phi \in \mathcal{F} \text{ s.t.} \\ (\mathcal{S}_P \phi)_- \in \mathcal{L}^1(Q)}} \left| \int \mathcal{S}_P \phi(x) \, dQ(x) \right|$$

where  $\mathcal{S}_P \phi(x) := (\nabla \log p)(x) \cdot \phi(x) + (\nabla \cdot \phi)(x)$ , where  $p(x) \propto q_0(x)p(y|x)$  is a density for  $P$ .

**$\therefore$  Stein discrepancy is measuring the size of a variational gradient**

## Computing gradient discrepancy

**Kernel gradient discrepancy:** Let  $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  be a kernel with associated Hilbert space  $\mathcal{H}_k$ . The Kernel Gradient Discrepancy (KGD) is defined as

$$\text{KGD}_k(Q) := \sup_{\substack{\|\phi\|_{\mathcal{H}_k} \leq 1 \text{ s.t.} \\ (\mathcal{T}_Q \phi)_- \in \mathcal{L}^1(Q)}} \left| \int \mathcal{T}_Q \phi(x) \, dQ(x) \right|.$$

$\therefore$  **KGD generalises kernel Stein discrepancy (KSD) to nonlinear  $\mathcal{L}$**

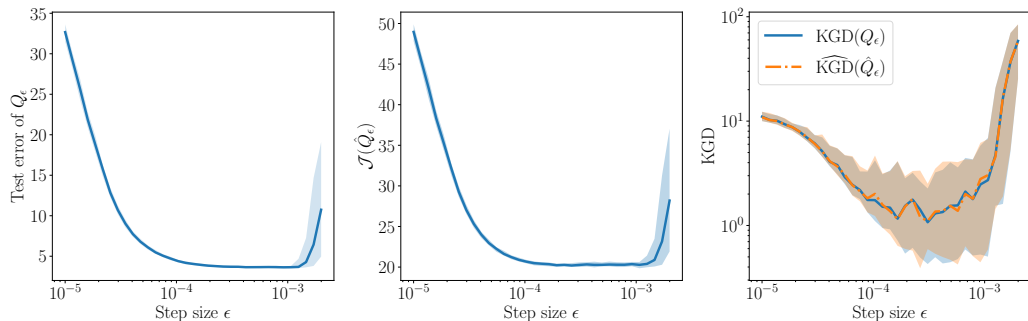
**Computable form of KGD:** Let  $s(Q)(\theta) := (\nabla \log q_0)(\theta) - \nabla_{\mathbf{V}} \mathcal{L}(Q)(\theta)$ . Then

$$\text{KGD}_k(Q) = \left( \iint k_Q(x, x') \, dQ(x) dQ(x') \right)^{1/2}$$

with the  $Q$ -dependent kernel

$$\begin{aligned} k_Q(x, x') &:= \nabla_1 \cdot \nabla_2 k(x, x') + \nabla_1 k(x, x') \cdot s(Q)(x') \\ &\quad + \nabla_2 k(x, x') \cdot s(Q)(x) + k(x, x') s(Q)(x) \cdot s(Q)(x'). \end{aligned}$$

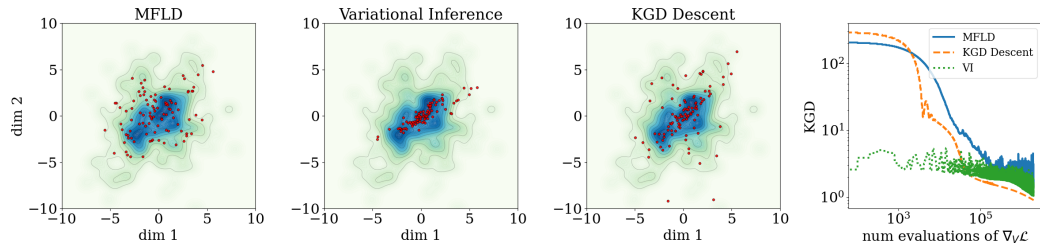
# Measuring approximation quality via KGD



**Figure:** Selecting the step size  $\epsilon$  in mean field Langevin dynamics for training a mean field neural network.

**Rather than use MFLD, why not directly optimise KGD?**

# New Algorithms Based on KGD (I)



**Figure:** Comparing MFLD with new algorithms based on KGD, in the setting of training mean field neural networks.

# Variational Gradient Descent

**Steepest descent:** Following Wang and Liu [2019], pick the vector field  $\phi_Q$  from a vector-valued reproducing kernel Hilbert space  $\mathcal{H}_k^d$  corresponding to steepest descent:

$$\phi_Q(\cdot) \propto \int \{k(x, \cdot)(\nabla \log q_0 - \nabla_V \mathcal{L}(Q))(x) + \nabla_1 k(x, \cdot)\} dQ(x),$$

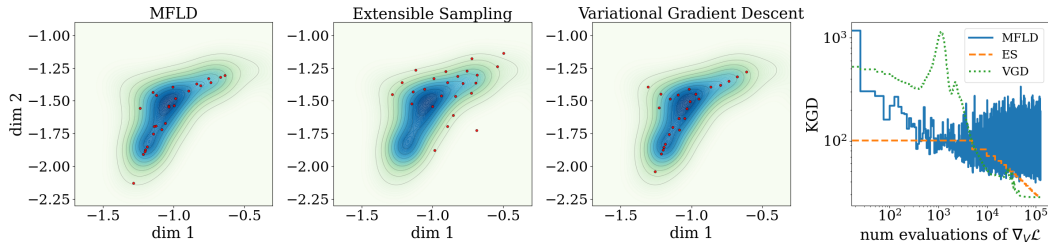
**Variational gradient descent:** Initialise  $\{x_i^0\}_{i=1}^N$  as independent samples from  $\mu_0$  at time  $t = 0$  and then update  $\{x_i^t\}_{i=1}^N$  deterministically, via the coupled system of ODEs

$$\frac{dx_k^t}{dt} = \frac{1}{N} \sum_{j=1}^N k(x_i^t, x_j^t)(\nabla \log q_0 - \nabla_V \mathcal{L}(Q_N^t))(x_j^t) + \nabla_1 k(x_j^t, x_i^t), \quad Q_N^t := \frac{1}{N} \sum_{j=1}^N \delta_{x_j^t}$$

up to a time horizon  $T$ .



## New Algorithms Based on KGD (II)



**Figure:** Comparing MFLD with new algorithms based on KGD, in the setting of prediction-centric uncertainty quantification.

# Properties of Kernel Gradient Discrepancy

## Informal Theorem (Properties of KGD)

*One can pick the kernel  $k$  such that the following hold:*

- ▶ **Identification:**  $\text{KGD}_k(Q) = 0$  iff  $Q$  a stationary point of  $\mathcal{J}$   
(roughly: convexity of  $\mathcal{L}$  implies a unique stationary point  $P$  of  $\mathcal{J}$ )
- ▶ **Continuity:**  $Q_n \xrightarrow{\alpha} Q$  implies  $\text{KGD}_k(Q_n) \rightarrow \text{KGD}_k(Q)$   
( $Q_n \xrightarrow{\alpha} Q$  means  $\int f \, dQ_n \rightarrow \int f \, dQ$  for all  $f(\theta) \lesssim 1 + \|\theta\|^\alpha$ )
- ▶ **Convergence control:**  $\text{KGD}_k(Q_n) \rightarrow 0$  implies  $Q_n \xrightarrow{\alpha} P$   
(requires a dissipativity condition on  $P$ )

**$\therefore$  minimisation of KGD leads to consistent approximation of the target**

# Summary

In a nutshell:

- ▶ (kernel) gradient discrepancy (KGD) enables approximation quality to be measured...
- ▶ ... and unlocks new classes of algorithms for  $\arg \min \mathcal{J}$
- ▶ can be considered a **nonlinear** generalisation of kernel Stein discrepancy (KSD)
- ▶ sheds light on KSD as measuring the size of a variational gradient

Open questions:

- ▶ All the usual challenges apply, e.g. high-dimensions, manifold-constrained targets, computational efficiency, mode collapse, etc.
- ▶ Stein's identity generalizes to the diffusion Stein operator [Gorham et al., 2019] – what would be the corresponding nonlinear generalization?

**Thank you for your attention!**

# References I

- P. G. Bissiri, C. C. Holmes, and S. G. Walker. A general framework for updating belief distributions. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 78(5):1103–1130, 2016.
- J. Gorham and L. Mackey. Measuring sample quality with Stein's method. *Advances in Neural Information Processing Systems*, 28, 2015.
- J. Gorham, A. B. Duncan, S. J. Vollmer, and L. Mackey. Measuring sample quality with diffusions. *The Annals of Applied Probability*, 29(5):2884–2928, 2019.
- Y. McLatchie, B.-E. Cherief-Abdellatif, D. T. Frazier, and J. Knoblauch. Predictively oriented posteriors. *arXiv preprint arXiv:2510.01915*, 2025.
- Z. Shen, J. Knoblauch, S. Power, and C. J. Oates. Prediction-centric uncertainty quantification via MMD. In *International Conference on Artificial Intelligence and Statistics*, pages 649–657. PMLR, 2025.
- D. Wang and Q. Liu. Nonlinear Stein variational gradient descent for learning diversified mixture models. In *International Conference on Machine Learning*, pages 6576–6585. PMLR, 2019.