

A Rigorous Link between Deep Ensembles and (Variational) Bayesian Methods

Veit D. Wild



UNIVERSITY OF
OXFORD

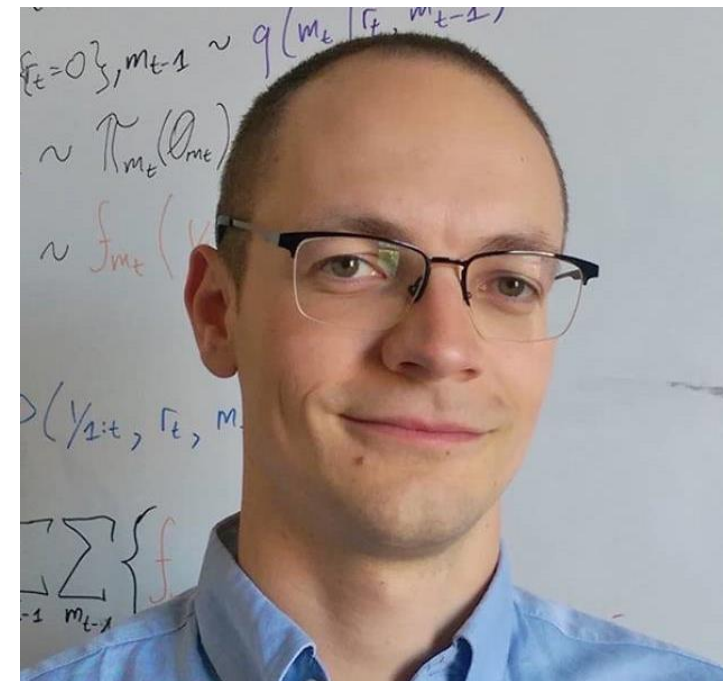
Department of Statistics



Sahra Ghalebikesabi



Dino Sejdinovic



Jeremias Knoblauch



Generalised Variational Inference

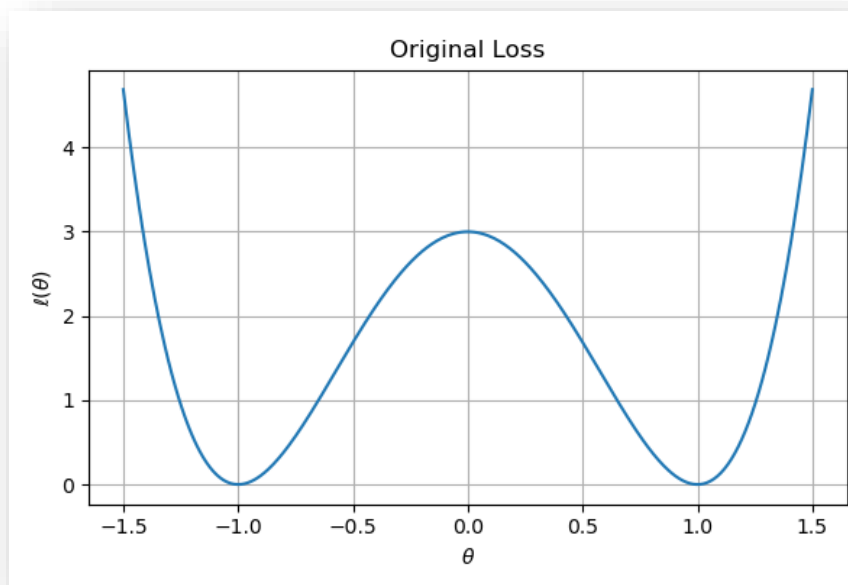
P prior/reference prob. measure
 $\lambda > 0$ reg. parameter
 $D(\cdot, P)$ convex regulariser

$$\theta^* \in \arg \min_{\theta \in \mathbb{R}^J} \ell(\theta)$$

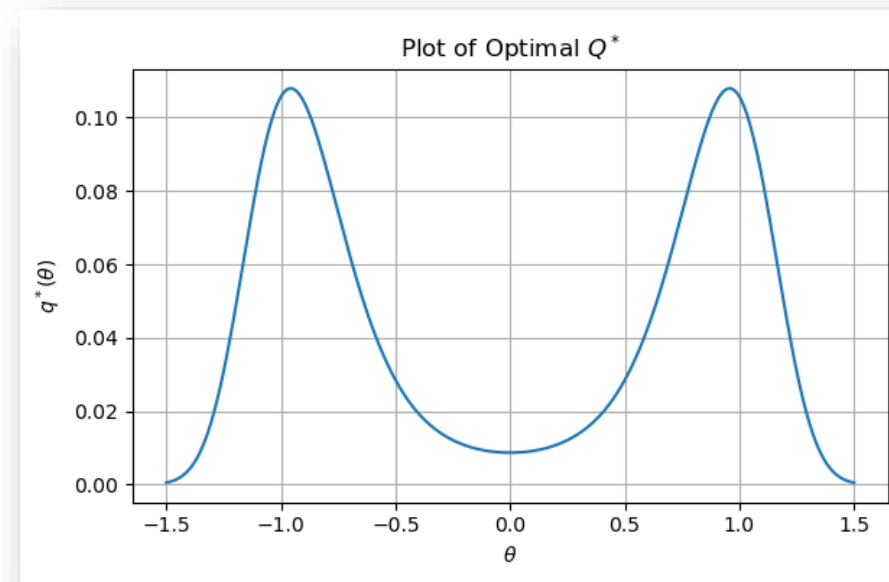
prob. lifting + convexify

$$L(Q) = \int \ell(\theta) dQ(\theta) + \lambda D(Q, P)$$

$$Q^* = \arg \min_{Q \in \mathcal{P}(\mathbb{R}^J)} L(Q)$$



choose λ, D, P



How to find the global minimiser?

Parameterised/Finite Dim GVI:

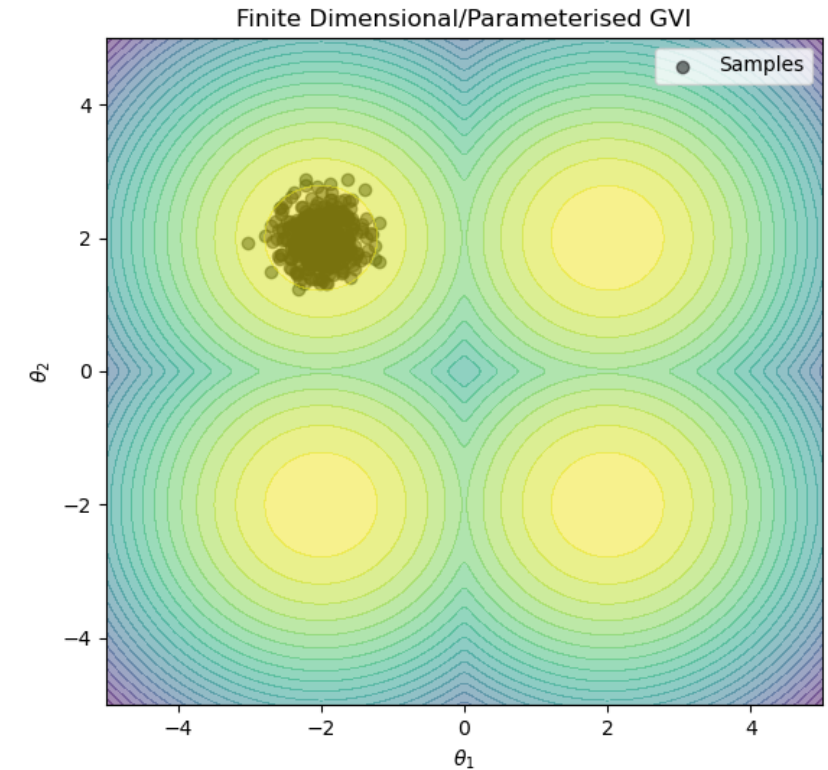
$$\mathcal{Q} := \{Q_\nu : \nu \in \Gamma\} \subset \mathcal{P}(\mathbb{R}^J)$$

e.g. : $Q_\nu = \mathcal{N}(\mu, \Sigma), \nu = (\mu, \Sigma)$

$$\tilde{L}(\nu) := L(Q_\nu) \longrightarrow \arg \min_{\nu \in \Gamma} \tilde{L}(\nu)$$

Alternative:

→ Gradient descent in infinite dimensions



Gradient Descent in Infinite Dimensions

Finite Dimensions

$$\theta_0 \in \mathbb{R}^J$$

$$\theta_{k+1} = \arg \min_{\theta \in \mathbb{R}^J} \left\{ \ell(\theta) + \frac{1}{2\eta} \|\theta - \theta_k\|_2^2 \right\}$$

$$\downarrow \quad \eta \rightarrow 0$$

$$\theta'(t) = -\nabla \ell(\theta(t))$$

Infinite Dimensions

$$Q_0 \in \mathcal{P}_2(\mathbb{R}^J)$$

$$Q_{k+1} := \arg \min_{Q \in \mathcal{P}_2(\mathbb{R}^J)} \left\{ L(Q) + \frac{1}{2\eta} \underbrace{W_2(Q, Q_k)^2}_{\text{Wasserstein metric}} \right\}$$


$$\downarrow \quad \eta \rightarrow 0$$

$$\partial_t q(t, \theta) = \nabla \cdot \left(q(t, \theta) \underbrace{\nabla_w L[Q(t)](\theta)}_{\text{Wasserstein gradient}} \right)$$

Hope: $q(t, \cdot) \rightarrow q^*$ for $t \rightarrow \infty$



Follow the WGF


$$L^{\text{fe}}(Q) := \int V(\theta) dQ(\theta) + \frac{\lambda_1}{2} \iint \kappa(\theta, \theta') dQ(\theta) dQ(\theta') + \lambda_2 \int \log q(\theta) q(\theta) d\theta$$

$$\nabla_W L^{\text{fe}}[Q](\theta) = \nabla V(\theta) + \lambda_1 \int (\nabla_1 \kappa)(\theta, \theta') dQ(\theta') + \lambda_2 \nabla \log q(\theta),$$

Step 1: Sample $N_E \in \mathbb{N}$ particles $\theta_1(0), \dots, \theta_{N_E}(0)$ independently from $Q_0 \in \mathcal{P}_2(\mathbb{R}^J)$.

Step 2: Evolve the particle θ_n by following the stochastic differential equation (SDE)

$$d\theta_n(t) = -\left(\nabla V(\theta_n(t)) + \frac{\lambda_1}{N_E} \sum_{j=1}^{N_E} (\nabla_1 \kappa)(\theta_n(t), \theta_j(t))\right) dt + \sqrt{2\lambda_2} dB_n(t),$$

for $n = 1, \dots, N_E$, and $\{B_n(t)\}_{t>0}$ stochastically independent Brownian motions.

Theoretical Analysis

$$L(Q) = \int \ell(\theta) dQ(\theta) + \lambda D(Q, P)$$



$$D(Q, P)$$

Step 1: Initialise $N_E \in \mathbb{N}$ particles $\theta_{1,0}, \dots, \theta_{N_E,0}$ from a user chosen initial distribution Q_0 .

Step 2: Evolve the particles forward in time according to

$$\theta_{n,k+1} = \theta_{n,k} - \eta \left(\nabla V(\theta_{n,k}) + \frac{\lambda_1}{N_E} \sum_{j=1}^{N_E} (\nabla_1 \kappa)(\theta_{n,k}, \theta_{j,k}) \right) + \sqrt{2\eta\lambda_2} Z_{n,k}$$

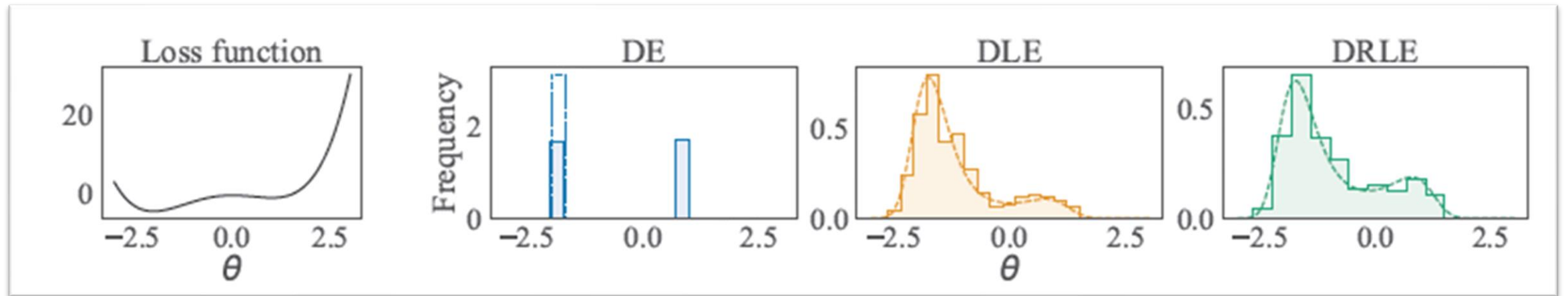
for $n = 1, \dots, N_E, k = 0, \dots, T - 1$ with $Z_{n,k} \sim \mathcal{N}(0, I_{J \times J})$.

$D(Q, P)$	$V(\theta)$	λ_1	λ_2	Method	Convergence
0	$\ell(\theta)$	0	0	DE	
KL(Q, P)	$\ell(\theta) - \lambda \log p(\theta)$	0	λ	DLE	
$\lambda \text{MMD}^2(Q, P) + \lambda' \text{KL}(Q, P)$	$\ell(\theta) - \underbrace{2\lambda\mu_P(\theta)}_{\text{kernel mean embedding of P}} - \lambda' \log p(\theta)$	2λ	λ'	DRLE	

kernel mean embedding of P



Visualisation



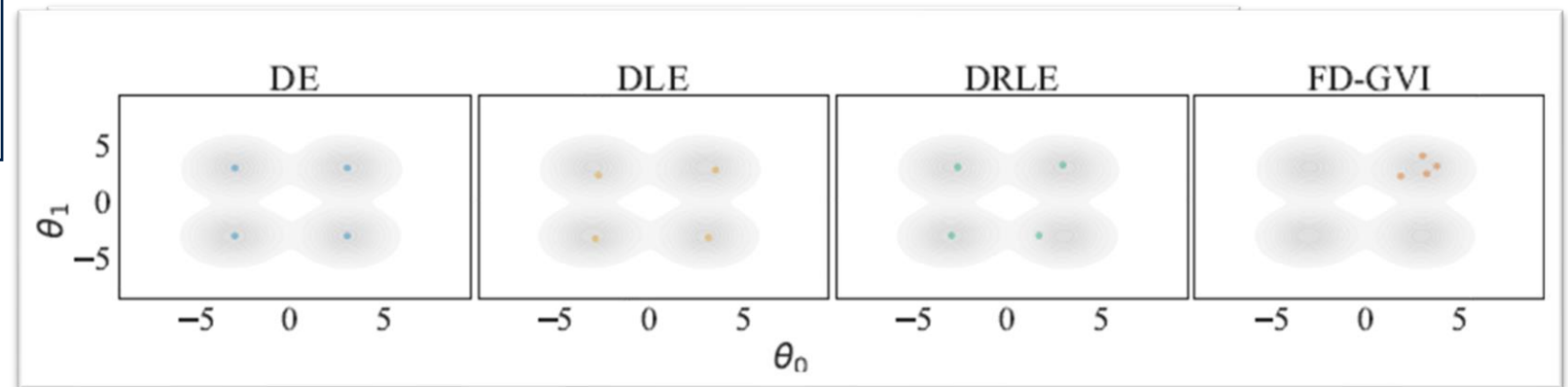
Theory meets Practice

	KIN8NM	CONCRETE	ENERGY	NAVAL	POWER	PROTEIN	WINE	YACHT
DE	$0.33_{\pm 0.1}$	$6.10_{\pm 0.3}$	$2.83_{\pm 0.2}$	$-0.40_{\pm 0.3}$	$13.70_{\pm 2.6}$	$11.22_{\pm 2.2}$	$14.65_{\pm 1.9}$	$2.20_{\pm 0.4}$
DLE	$13.25_{\pm 4.3}$	$5.11_{\pm 0.2}$	$2.43_{\pm 0.1}$	$3.46_{\pm 2.4}$	$13.87_{\pm 2.3}$	$43.20_{\pm 12.5}$	$13.73_{\pm 1.4}$	$1.64_{\pm 0.1}$
DRLE	$0.46_{\pm 0.1}$	$8.30_{\pm 0.6}$	$4.01_{\pm 0.3}$	$-3.04_{\pm 0.2}$	$23.21_{\pm 2.0}$	$48.80_{\pm 2.1}$	$7.13_{\pm 0.6}$	$7.80_{\pm 2.7}$

No difference, but why?

- Finite time?
- Discretisation error?
- Batch estimation?

nr. sample size \ll nr. modes



Thank you for your attention!



Deep Ensembles

$$d\theta_n(t) = -\left(\nabla V(\theta_n(t)) + \frac{\lambda_1}{N_E} \sum_{j=1}^{N_E} (\nabla_1 \kappa)(\theta_n(t), \theta_j(t))\right) dt + \sqrt{2\lambda_2} dB_n(t)$$

$$L(Q) = \int \ell(\theta) dQ(\theta) + \lambda D(Q, P) \xrightarrow{\text{choose}} D(Q, P) = 0$$

$V(\theta)$	λ_1	λ_2
$\ell(\theta)$	0	0

Theorem 1. *If ℓ has countably many local minima $\{m_i : i \in \mathbb{N}\}$, then it holds independently for each $n = 1, \dots, N_E$ that*

$$\theta_n(t) \xrightarrow{\mathcal{D}} \sum_{i=1}^{\infty} Q_0(\Theta_i) \delta_{m_i} =: Q_{\infty}$$

for $t \rightarrow \infty$. Here $\xrightarrow{\mathcal{D}}$ denotes convergence in distribution and $\Theta_i = \{\theta \in \mathbb{R}^J : \lim_{t \rightarrow \infty} \theta_*(t) = m_i \text{ and } \theta_*(0) = \theta\}$ denotes the domain of attraction for m_i with respect to the gradient flow θ_* .



Deep Langevin Ensembles

$$d\theta_n(t) = -\left(\nabla V(\theta_n(t)) + \frac{\lambda_1}{N_E} \sum_{j=1}^{N_E} (\nabla_1 \kappa)(\theta_n(t), \theta_j(t))\right) dt + \sqrt{2\lambda_2} dB_n(t)$$

$$L(Q) = \int \ell(\theta) dQ(\theta) + \lambda D(Q, P)$$

choose $\longrightarrow D(Q, P) = \text{KL}(Q, P)$

$V(\theta)$	λ_1	λ_2
$\ell(\theta) - \lambda \log p(\theta)$	0	λ

Theorem 2. For each $n = 1, \dots, N_E$ independently

$$\theta_n(t) \xrightarrow{\mathcal{D}} Q^*$$

with $q^*(\theta) \propto \exp\left(-\frac{1}{\lambda}\ell(\theta)\right)p(\theta)$.



Deep Rep. Langevin Ensembles

$$d\theta_n(t) = -\left(\nabla V(\theta_n(t)) + \frac{\lambda_1}{N_E} \sum_{j=1}^{N_E} (\nabla_1 \kappa)(\theta_n(t), \theta_j(t))\right) dt + \sqrt{2\lambda_2} dB_n(t)$$

$$L(Q) = \int \ell(\theta) dQ(\theta) + D(Q, P)$$

choose
 $\longrightarrow D(Q, P) = \lambda \text{MMD}^2(Q, P) + \lambda' \text{KL}(Q, P)$

$V(\theta)$	λ_1	λ_2
$\ell(\theta) - 2\lambda\mu_P(\theta) - \lambda' \log p(\theta)$	2λ	λ'

Theorem 3. Let $Q^{n, N_E}(t)$ be the distribution of $\theta_n(t)$, $n = 1, \dots, N_E$, generated via the WGF. Then

$$Q^{n, N_E}(t) \xrightarrow{\mathcal{D}} Q^*$$

for each $n = 1, \dots, N_E$ and as $N_E \rightarrow \infty$, $t \rightarrow \infty$.

