

511 **A Existence and uniqueness of global minimiser**

512 In this section, we discuss assumptions under which the global minimiser of the optimisation problem
513 lem

$$L(Q) = \int \ell(\theta) dQ(\theta) + \lambda D(Q, P) \quad (8)$$

514 over $\mathcal{P}(\mathbb{R}^J)$ exists and is unique. We assume throughout that the optimisation problem is not patho-
515 logical, in the sense that there exists a measure $\widehat{Q} \in \mathcal{P}(\mathbb{R}^J)$ such that $L(\widehat{Q}) < \infty$. This is in
516 applications often trivial to verify. A good candidate for \widehat{Q} is typically the reference measure P .

517 **Loss assumptions** Let $\ell : \mathbb{R}^J \rightarrow \mathbb{R}$ be a loss satisfying the following assumptions:

518 (L1) The loss ℓ is bounded from below which means that

$$c := \inf \{ \ell(\theta) : \theta \in \mathbb{R}^J \} > -\infty. \quad (9)$$

519 (L2) The loss is norm-coercive which means that

$$\ell(\theta) \rightarrow \infty \quad (10)$$

520 if $\|\theta\| \rightarrow \infty$.

521 (L3) The loss ℓ is lower semi-continuous which means that

$$\liminf_{\theta \rightarrow \theta_0} \ell(\theta) \geq \ell(\theta_0) \quad (11)$$

522 for all $\theta_0 \in \mathbb{R}^J$.

523 **Regulariser assumptions** Let $D : \mathcal{P}(\mathbb{R}^J) \times \mathcal{P}(\mathbb{R}^J) \rightarrow [0, \infty]$ be a regulariser and $P \in \mathcal{P}(\mathbb{R}^J)$
524 a reference measure. We define $D_P(\cdot) := D(\cdot, P)$ for notational convenience. We assume the
525 following for D_P :

526 (D1) The function D_P is lower semi-continuous w.r.t. to the topology of weak-convergence, i.e.
527 for all sequences $(Q_n)_{n \in \mathbb{N}} \subset \mathcal{P}(\mathbb{R}^J)$ and all Q with $D_P(Q) < \infty$, it holds that $Q_n \xrightarrow{\mathcal{D}} Q$
528 implies

$$\liminf_{n \rightarrow \infty} D_P(Q_n) \geq D_P(Q). \quad (12)$$

529 Here, $\xrightarrow{\mathcal{D}}$ denotes convergence in distribution.

530 (D2) D_P is strictly convex, i.e. for all $Q_1 \neq Q_2 \in \mathcal{P}(\mathbb{R}^J)$ with $D_P(Q_1) < \infty$ and $D_P(Q_2) < \infty$, it holds that

$$D_P(\alpha Q_1 + (1 - \alpha) Q_2) < \alpha D_P(Q_1) + (1 - \alpha) D_P(Q_2) \quad (13)$$

532 with $\alpha \in (0, 1)$.

533 The next theorem provides an existence result for the optimisation problem (8). The result is similar
534 in spirit to Lemma 2.1 in [Knoblauch \(2021\)](#) with the important difference that our assumptions are
535 easier to verify, since they are formulated in terms of ℓ and D_P .

536 **Theorem 3** (Existence of global minimiser). *Under the assumptions (L1)-(L3) and (D1) there exists*
537 *a probability measure $Q^* \in \mathcal{P}(\mathbb{R}^J)$ with*

$$L(Q^*) = \inf \{ L(Q) : Q \in \mathcal{P}(\mathbb{R}^J) \}. \quad (14)$$

538 *Proof.* Let $c > -\infty$ be the lower bound for ℓ . It follows immediately that $L(Q) \geq c$ for all
539 $Q \in \mathcal{P}(\mathbb{R}^J)$ since $D(P, Q) \geq 0$. As a consequence we know that

$$\infty > L^* := \inf \{ L(Q) : Q \in \mathcal{P}(\mathbb{R}^J) \} \geq c > -\infty. \quad (15)$$

540 By definition of the infimum we can construct a sequence $l_n = L(Q_n) \in \mathbb{R}$ in the image of L such

$$l_n \rightarrow L^* \quad (16)$$

541 for $n \rightarrow \infty$. We now show by contradiction that the corresponding sequence $(Q_n) \subset \mathcal{P}(\mathbb{R}^J)$ is
542 *tight*³. Assume that (Q_n) is not tight. By definition we can then find an $\epsilon > 0$ such that for each
543 $k \in \mathbb{N}$ there exists $n = n_k \in \mathbb{N}$ with $Q_{n_k}([-k, k]^J) \leq 1 - \epsilon$. We set $A_k := [-k, k]^J \subset \mathbb{R}^J$ and
544 obtain

$$l_{n_k} = L(Q_{n_k}) \tag{17}$$

$$= \int_{A_k} \ell(\theta) dQ_{n_k}(\theta) + \int_{\mathbb{R}^J \setminus A_k} \ell(\theta) dQ_{n_k}(\theta) + \lambda D(Q, P) \tag{18}$$

$$\geq \int_{A_k} \ell(\theta) dQ_{n_k}(\theta) + \int_{\mathbb{R}^J \setminus A_k} \ell(\theta) dQ_{n_k}(\theta) \tag{19}$$

$$\geq cQ_{n_k}(A_k) + \inf \{ \ell(\theta) : \theta \in \mathbb{R}^J \setminus A_k \} Q_{n_k}(\mathbb{R}^J \setminus A_k) \tag{20}$$

$$\geq cQ_{n_k}(A_k) + \epsilon \inf \{ \ell(\theta) : \theta \in \mathbb{R}^J \setminus A_k \}. \tag{21}$$

545 Due to the coerciveness of ℓ , we know that $\inf \{ \ell(\theta) : \theta \in \mathbb{R}^J \setminus A_k \} \rightarrow \infty$ for $k \rightarrow \infty$ and there-
546 fore $l_{n_k} \rightarrow \infty$ for $k \rightarrow \infty$. However, this is a contradiction: The sequence (l_n) is convergent and
547 therefore in particular bounded. As a consequence, it cannot contain the unbounded sub-sequence
548 (l_{n_k}) . It follows that the sequence (Q_n) is tight. By Prokhorov's theorem we can now extract a sub
549 sequence (Q_{n_k}) of (Q_n) and a measure $Q^* \in \mathcal{P}(\mathbb{R}^J)$ such that

$$Q_{n_k} \xrightarrow{\mathcal{D}} Q^* \tag{22}$$

550 for $k \rightarrow \infty$. Due to Lemma 5.1.7 in [Ambrosio et al. \(2005\)](#) the lower semi-continuity of ℓ implies
551 that $Q \mapsto \int \ell(\theta) dQ(\theta)$ is lower semi-continuous. This combined with the lower semi-continuity of
552 D_P gives

$$\liminf_{k \rightarrow \infty} L(Q_{n_k}) \geq L(Q^*). \tag{23}$$

553 From this it immediately follows that

$$L(Q^*) \leq \liminf_{k \rightarrow \infty} L(Q_{n_k}) = L^*, \tag{24}$$

554 but by definition L^* is the global minimum of L which implies $L^* \leq L(Q^*)$. We therefore conclude
555 that $L(Q^*) = L^*$. \square

556 Theorem [3](#) only shows the existence of a global minimiser. In order to show uniqueness we use
557 the convexity assumption (D2). The proof is the same as in finite dimensions and only included for
558 completeness.

559 **Theorem 4** (Uniqueness of global minimiser). *Assume that (D2) holds. Then, the global minimiser*
560 *of L is unique (whenever it exists).*

561 *Proof.* Assume there exists two probability measures $Q_1, Q_2 \in \mathcal{P}(\mathbb{R}^J)$ such that

$$L(Q_1) = L^* = L(Q_2). \tag{25}$$

562 where $\infty > L^* := \inf \{ L(Q) : Q \in \mathcal{P}(\mathbb{R}^J) \} > -\infty$. We define the probability measure
563 $Q_3 := \frac{1}{2}Q_1 + \frac{1}{2}Q_2$. By strict convexity we obtain

$$L(Q_3) < \frac{1}{2}L(Q_1) + \frac{1}{2}L(Q_2) = L^*, \tag{26}$$

564 which is a contradiction to Q_1 and Q_2 being global minimisers. \square

565 Note that in the literature on GVI ([Knoblauch et al., 2022](#)) it is common to assume that the regulariser
566 is definite, i.e.

$$D(P, Q) = 0 \iff P = Q \tag{27}$$

567 for all $P, Q \in \mathcal{P}(\mathbb{R}^J)$. We did not use this assumption in neither Theorem [3](#) nor Theorem [4](#).
568 However, the next lemma shows that it is basically implied by strict convexity.

³A sequence of probability measures (Q_n) is called tight if and only if for every $\epsilon > 0$ there exists a compact set $K \in \mathbb{R}^J$ such that for all $n \in \mathbb{N}$ holds: $Q_n(K) > 1 - \epsilon$.

569 **Lemma 1.** Let $D_P : \mathcal{P}(\mathbb{R}^J) \rightarrow [0, \infty]$ be strictly convex and assume further $D(Q, Q) =$
570 0 for all $Q \in \mathcal{P}(\mathbb{R}^J)$. Then it follows that $D(Q, P) = 0$ implies $P = Q$.

571 *Proof.* We prove the claim by contradiction. Assume that there exists $P \neq Q$ such that $D(P, Q) =$
572 0 . The strict convexity and $D(P, P) = 0$ imply combined that

$$D\left(\frac{1}{2}P + \frac{1}{2}Q, P\right) < \frac{1}{2}D(P, P) + \frac{1}{2}D(Q, P) \quad (28)$$

$$= 0. \quad (29)$$

573 However, we know that $D(\frac{1}{2}P + \frac{1}{2}Q, P) \geq 0$ by assumption. This is a contradiction. \square

574 **Discussion on loss assumptions** The assumptions on the loss ℓ in (L1) and (L3) are rather weak.
575 Typically loss functions in machine learning are bounded from below and continuous (and there-
576 fore in particular lower semi-continuous). However, norm-coercivity can be violated. Consider for
577 example the squared loss

$$\ell(\theta) := \sum_{n=1}^N (y_n - f_\theta(x_n))^2, \quad (30)$$

578 where f_θ is the parametrisation of a neural network with one hidden layer, i.e. $\theta = (w, A)$ and

$$f_\theta(x) = w^T \sigma(Ax), \quad (31)$$

579 where $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is an activation function which is applied pointwise to the vector Ax and has the
580 property that $\sigma(0) = 0$. It is now possible to find a sequence of parameters $(\theta_k)_{k \in \mathbb{N}} \subset \mathbb{R}^J$ with
581 $\|\theta_k\| \rightarrow \infty$ such that $\ell(\theta_k)$ does not converge to infinity. Define $w_k := k(1 \dots 1)$, $A_k := 0$ and
582 $\theta_k = (w_k, A_k)$ for $k \in \mathbb{N}$. Then we obviously have that

$$\|\theta_k\| = \|w_k\| \rightarrow \infty \quad (32)$$

583 for $k \rightarrow \infty$ but

$$\ell(\theta_k) = \sum_{n=1}^N (y_n - f_{\theta_k}(x_n))^2 \quad (33)$$

$$= \sum_{n=1}^N (y_n - w_k^T \sigma(0))^2 \quad (34)$$

$$= \sum_{n=1}^N y_n^2, \quad (35)$$

584 which is constant and therefore does not converge to ∞ . A similar, but notationally more involved,
585 construction can be made for neural networks with more than one hidden layer. However, this is an
586 issue that can be easily resolved by adding what is known as weight decay to the loss. For example,
587 consider for $\gamma > 0$ the loss

$$\ell(\theta) := \sum_{n=1}^N (y_n - f_\theta(x_n))^2 + \gamma \|\theta\|^2 \quad (36)$$

588 with weight decay. This loss is by construction norm-coercive and therefore the previous existence
589 proof applies.

590 **Discussion on regulariser assumptions** The assumptions (D1) and (D2) are quite weak. The KL-
591 divergence for example is known to be lower semi-continuous (Polyanskiy and Wu, 2014, Theorem
592 3.7) and strictly convex (Polyanskiy and Wu, 2014, Theorem 4.1). This immediately implies lower
593 semi-continuity and convexity of $\text{KL}(\cdot, P)$ for any fixed P . The MMD is also known to be strictly
594 convex (Arbel et al., 2019, Lemma 25), whenever it is well-defined, which can be guaranteed under
595 weak assumptions on κ (Muandet et al., 2017, Lemma 3.1). The lower semi-continuity properties
596 also depend on the kernel κ . However, for bounded kernels it is trivial to verify. We include the
597 proof for completeness, but assume this has been shown before elsewhere.

598 **Lemma 2.** Let the kernel $\kappa : \mathbb{R}^J \times \mathbb{R}^J$ be continuous and bounded: $\|\kappa\|_\infty :=$
599 $\sup_{\theta, \theta' \in \mathbb{R}^J} |\kappa(\theta, \theta')| < \infty$ and P be fixed. Then $\text{MMD}(\cdot, P)$ is continuous and therefore, in partic-
600 ular, lower semi-continuous.

601 *Proof.* Let $(Q_n)_{n \in \mathbb{N}}$ and Q^* be such that

$$Q_n \xrightarrow{\mathcal{D}} Q^* \quad (37)$$

602 for $n \rightarrow \infty$. This immediately implies that

$$Q_n \otimes Q_n \xrightarrow{\mathcal{D}} Q^* \otimes Q^* \quad (38)$$

603 for $n \rightarrow \infty$, where $Q^* \otimes Q^*$ denotes the product measure of Q^* with itself. Further, note that
604 the kernel mean embedding μ_P is continuous as integral with respect to the second component of a
605 continuous function and bounded since

$$|\mu_P(\theta)| = \left| \int \kappa(\theta, \theta') dP(\theta') \right| \quad (39)$$

$$\leq \int |\kappa(\theta, \theta')| dP(\theta') \quad (40)$$

$$\leq \|\kappa\|_\infty. \quad (41)$$

606 By the definition of weak convergence for measures, we therefore have

$$\iint \kappa(\theta, \theta') d(Q_n \otimes Q_n)(\theta, \theta') \longrightarrow \iint \kappa(\theta, \theta') d(Q^* \otimes Q^*)(\theta, \theta') \quad (42)$$

$$\int \mu_P(\theta) dQ_n(\theta) \longrightarrow \int \mu_P(\theta) dQ^*(\theta) \quad (43)$$

607 for $n \rightarrow \infty$. This immediately implies continuity of $\text{MMD}(\cdot, P)$ with respect to the topology of
608 weak convergence. \square

609 Notice that most kernels common in machine learning, such as the squared exponential or the Matérn
610 kernel, are continuous and bounded and therefore Lemma 2 applies.

611 **Remark 1.** The astute reader may have noticed that our existence proof only guarantees the ex-
612 istence of measure $Q^* \in \mathcal{P}(\mathbb{R}^J)$. However, the Wasserstein gradient flow is by definition only
613 formulated in the space of probability measures with finite second moment, denoted $\mathcal{P}_2(\mathbb{R}^J)$. As-
614 sumptions which guarantee that $Q^* \in \mathcal{P}_2(\mathbb{R}^J)$ are easy to formulate. For example, we can require
615 that there exists $C > 0$ and $R > 0$ such that the loss ℓ satisfies

$$|\ell(\theta)| > C\|\theta\|^2 \quad (44)$$

616 for all $\|\theta\| > R$. This immediately implies that $Q^* \in \mathcal{P}_2(\mathbb{R}^J)$ since otherwise

$$\int |\ell(\theta)| dQ^*(\theta) = \infty \quad (45)$$

617 gives a contradiction to the finiteness of $L(Q^*)$. However, even if (44) is violated, the reference
618 measure P may still guarantee that $Q^* \in \mathcal{P}_2(\mathbb{R}^J)$. For example, if $P \in \mathcal{P}_2(\mathbb{R}^J)$, then $D_P(Q^*)$ will
619 typically be large if $Q^* \notin \mathcal{P}_2(\mathbb{R}^J)$ and the global minimiser is therefore in a sense *unlikely* to have
620 fat tails. We therefore assume $Q^* \in \mathcal{P}_2(\mathbb{R}^J)$ throughout the paper and consider it to be a minor
621 practical concern.

622 B Realising the Wasserstein gradient flow

623 In this section, we identify a suitable stochastic process that allows us to follow the WGF.

624 Let $L^{\text{fe}} : \mathcal{P}(\mathbb{R}^J) \rightarrow (-\infty, \infty]$ be the free energy discussed in Section 3.2 given as

$$L^{\text{fe}}(Q) := \int V(\theta) dQ(\theta) + \frac{\lambda_1}{2} \int \kappa(\theta, \theta') dQ(\theta) dQ(\theta') + \lambda_2 \int \log(q(\theta)) q(\theta) d\theta, \quad (46)$$

625 where $\lambda_1, \lambda_2 \geq 0$ are constants, $V : \mathbb{R}^J \rightarrow \mathbb{R}$ is the potential, $\kappa : \mathbb{R}^J \times \mathbb{R}^J \rightarrow \mathbb{R}$ is symmetric. We
626 will write L for L^{fe} from now on to simplify notation. The Wasserstein gradient of L is given as (cf.
627 Chapter 9.1 [Villani, 2003](#) Equation 9.4)

$$\nabla_W L[Q](\theta) = \nabla V(\theta) + \lambda_1(\nabla_1 \kappa * Q)(\theta) + \lambda_2 \nabla \log(q(\theta)), \quad (47)$$

628 where $\nabla_1 \kappa : \mathbb{R}^J \times \mathbb{R}^J \rightarrow \mathbb{R}^J$ is the (vector-valued) derivative of κ with respect to the first compo-
629 nent, ∇ denotes the euclidean gradient with respect to θ and $(\nabla_1 \kappa * Q)(\theta) := \int \nabla_1 \kappa(\theta, \theta') dQ(\theta')$
630 for $\theta \in \mathbb{R}^J$. The corresponding Wasserstein gradient flow is therefore given as (cf. Chapter 9.1
631 [Villani, 2003](#) Equation 9.3)

$$\partial_t q(t, \theta) = \nabla \cdot \left(q(t, \theta) (\nabla V(\theta) + \lambda_1(\nabla_1 \kappa * Q)(\theta) + \lambda_2 \nabla \log(q_t(\theta))) \right). \quad (48)$$

632 In general the probability density evolution of a stochastic process is—via the Fokker-Planck
633 equation—associated with the adjoint of the (infinitesimal) generator of the stochastic process. We
634 will therefore try to identify the generator associated to the density evolution in [\(48\)](#). To this end
635 let $h \in C_c^2(\mathbb{R}^J, \mathbb{R})$ where $C_c^2(\mathbb{R}^J, \mathbb{R})$ denotes the space of twice continuously differentiable func-
636 tions with compact support. We multiply both sides of [\(48\)](#) with h , integrate, and apply the partial
637 integration rule to obtain

$$\frac{d}{dt} \int h(\theta) q(t, \theta) d\theta = - \int \nabla_W L[Q(t)](\theta) \cdot \nabla h(\theta) q(t, \theta) d\theta. \quad (49)$$

$$= - \int (\nabla V(\theta) + \lambda_1(\nabla_1 \kappa * Q_t)(\theta)) \cdot \nabla h(\theta) dQ_t(\theta) \quad (50)$$

$$- \lambda_2 \int \nabla \log(q_t(\theta)) \cdot \nabla h(\theta) dQ_t(\theta). \quad (51)$$

638 By chain-rule and partial integration, [\(51\)](#) can be rewritten as

$$- \lambda_2 \int \nabla \log(q_t(\theta)) \cdot \nabla h(\theta) dQ_t(\theta) = - \lambda_2 \int \nabla q_t(\theta) \cdot \nabla h(\theta) d\theta \quad (52)$$

$$= \lambda_2 \int \Delta h(\theta) dQ_t(\theta). \quad (53)$$

639 Putting everything together, we obtain

$$\frac{d}{dt} \int h(\theta) q(t, \theta) d\theta = \int (A[Q(t)]h)(\theta) dQ_t(\theta), \quad (54)$$

640 where $\{A[Q]\}_{Q \in \mathcal{P}(\mathbb{R}^J)}$ is a family of operators defined as

$$(A[Q]h)(\theta) := - \left(\nabla V(\theta) + \lambda_1(\nabla_1 \kappa * Q)(\theta) \right) \cdot \nabla h(\theta) + \lambda_2 \Delta h. \quad (55)$$

641 for $h \in C_c^2(\mathbb{R}^J, \mathbb{R})$. The reader may recognize this operator family as the generator of a so called
642 *nonlinear Markov processes* ([Kolokoltsov, 2010](#), Chapter 1.4). The nonlinearity in this case refers
643 to the dependency on the measure Q . Linear Markov processes have no measure-dependency. This
644 family of generators corresponds to a McKean-Vlasov process of the form

$$d\theta(t) = - \left(\nabla V(\theta(t)) + \lambda_1(\nabla_1 \kappa * Q_t)(\theta(t)) \right) dt + \sqrt{2\lambda_2} dB(t), \quad (56)$$

645 where $(B(t))_{t>0}$ is a Brownian motion and Q_t the law of $\theta(t)$. In other words: The solution to
646 [\(56\)](#) has the time marginals $Q(t)$ such that [\(54\)](#) holds for every $h \in C_c^2(\mathbb{R}^J, \mathbb{R})$. Furthermore, the
647 corresponding pdfs $(q(t))$ satisfy the nonlinear Fokker-Planck equation given as

$$\partial_t q_t = A^*[Q_t]q_t, \quad (57)$$

648 where $A^*[Q]$ denotes the L^2 -adjoint of the operator $A[Q]$ and is given as

$$(A^*[Q]h)(\theta) = \nabla \cdot \left(h(\theta) (\nabla V(\theta) + \lambda_1(\nabla_1 \kappa * Q)(\theta) + \lambda_2 \nabla \log(h(\theta))) \right) \quad (58)$$

649 for $h \in C_c^2(\mathbb{R}^J, \mathbb{R})$ (Barbu and Röckner, 2020), cf. equation (1.1)-(1.4). Note that (57) corresponds
 650 exactly to the Wasserstein gradient flow equation in (48). We can therefore follow the WGF by
 651 simulating solutions to (56).

652 The standard approach to simulate solutions to (56) (Veretennikov, 2006) is to use an ensemble of
 653 interacting particles. Formally, we replace $Q(t)$ by $\frac{1}{N_E} \sum_{n=1}^{N_E} \delta_{\theta_n(t)}$ and obtain

$$d\theta_n(t) = -\left(\nabla V(\theta_n(t)) + \frac{\lambda_1}{N_E} \sum_{j=1}^{N_E} (\nabla_1 \kappa)(\theta_n(t), \theta_j(t))\right) dt + \sqrt{2\lambda_2} dB_n(t) \quad (59)$$

654 for $n = 1, \dots, N_E$ where $N_E \in \mathbb{N}$ denotes the number of particles. The Euler-Maruyama approxi-
 655 mation of (59) leads to the final algorithm:

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

657 **Step 2:** Evolve the particles forward in time according to

$$\theta_{n,k+1} = \theta_{n,k} - \eta \left(\nabla V(\theta_{n,k}) + \frac{\lambda_1}{N_E} \sum_{j=1}^{N_E} (\nabla_1 \kappa)(\theta_{n,k}, \theta_{j,k}) \right) + \sqrt{2\eta\lambda_2} Z_{n,k} \quad (60)$$

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

659 Note that $\theta_{n,k}$ is thought of as approximation of $\theta_n(t)$ at position $t = k\eta$. Furthermore, as dis-
 660 cussed in Section 4, various choices of V , λ_1 and λ_2 allow us to implement the WGF for different
 661 regularised optimisation problems in the space of probability measures. This is summarised below:

- 662 • Deep ensembles: $V(\theta) = \ell(\theta)$, $\lambda_1 = 0$, $\lambda_2 = 0$
- 663 • Deep Langevin ensembles: $V(\theta) = \ell(\theta) - \lambda \log p(\theta)$, $\lambda_1 := 0$, $\lambda := \lambda_2$
- 664 • Deep repulsive Langevin ensembles: $V(\theta) = \ell(\theta) - \lambda_1 \log p(\theta) - \lambda_2 \mu_P(\theta)$

665 C Asymptotic distribution of particles: unregularised objective

666 In this section, we investigate the asymptotic distribution of the WGF for the objective

$$L(Q) := \int \ell(\theta) dQ(\theta) \quad (61)$$

667 for $Q \in \mathcal{P}(\mathbb{R}^J)$. The associated particle method is:

- 668 • Sample $\theta_1(0), \dots, \theta_{N_E}(0)$ independently from Q_0 .
- 669 • Simulate (deterministically) $\theta'_n(t) = -\nabla \ell(\theta_n(t))$ for $n = 1, \dots, N_E$.

670 We start by introducing some notation for the deterministic gradient system. Let $\phi^t(\theta_0)$ denote the
 671 solution to the ordinary differential equation (ODE)

$$\theta(0) = \theta_0 \in \mathbb{R}^J \quad (62)$$

$$\theta'(t) = -\nabla \ell(\theta(t)) \quad (63)$$

672 at time $t > 0$. In a first step, we show the following lemma, which is a simple application of the
 673 famous Lojasiewicz theorem (Colding and Minicozzi II, 2014), and the fact that Lebesgue almost
 674 every initialisation leads to a local minimum (Lee et al., 2016).

675 **Lemma 3.** Assume $\ell : \mathbb{R}^J \rightarrow \mathbb{R}$ is norm-coercive and satisfies the Lojasiewicz inequality, i.e. for
 676 every $\theta \in \mathbb{R}^J$ exists an environment U of θ and constants $0 < \gamma < 1$ and $C > 0$ such that

$$|\ell(\theta) - \ell(\bar{\theta})|^\gamma < C |\nabla \ell(\theta)|. \quad (64)$$

677 for all $\bar{\theta} \in U$. Then we know that $\phi^t(\theta_0)$ converges for $t \rightarrow \infty$ to a local minimum of ℓ for Lebesgue
 678 almost every $\theta_0 \in \mathbb{R}^J$.

679 *Proof.* First we show that $t \mapsto \phi^t(\theta_0)$ is bounded. We proof this by contradiction. Assume that
 680 $\phi^t(\theta_0)$ is unbounded. Then there exists a subsequence $(t_n)_{n \in \mathbb{N}} \subset [0, \infty)$ with $t_n \rightarrow \infty$ for $n \rightarrow \infty$
 681 such that

$$|\phi^{t_n}(\theta_0)| \rightarrow \infty \quad (65)$$

682 for $n \rightarrow \infty$. The norm-coercivity immediately implies that

$$\ell(\phi^{t_n}(\theta_0)) \rightarrow \infty \quad (66)$$

683 for $n \rightarrow \infty$. However, this contradicts

$$\ell(\phi^t(\theta_0)) \leq \ell(\phi^0(\theta_0)) = \ell(\theta_0) < \infty, \quad (67)$$

684 where the first inequality follows from the fact that $t \mapsto \ell(\phi^t(\theta_0))$ is decreasing, which is a conse-
 685 quence of

$$\frac{d}{dt} \ell(\phi^t(\theta_0)) = \nabla \ell(\phi^t(\theta_0)) \frac{d}{dt} \phi^t(\theta_0) \quad (68)$$

$$= -|\nabla \ell(\phi^t(\theta_0))|^2 \leq 0. \quad (69)$$

686 Hence $t \mapsto \phi^t(\theta_0)$ is bounded. By the Bolzano-Weierstrass theorem we can find a sequence
 687 $(t_n)_{n \in \mathbb{N}} \subset [0, \infty)$ with $t_n \rightarrow \infty$ and a point $\theta_\infty \in \mathbb{R}^J$ such that

$$\phi^{t_n}(\theta_0) \rightarrow \theta_\infty \quad (70)$$

688 for $n \rightarrow \infty$. Hence $(\phi^t(\theta_0))_{t>0}$ has the accumulation point θ_∞ . The Lojasiewicz theorem (Colding
 689 and Minicozzi II, 2014) allows us to deduce that

$$\phi^t(\theta_0) \rightarrow \theta_\infty \quad (71)$$

690 for $t \rightarrow \infty$, and that θ_∞ satisfies $\nabla \ell(\theta_\infty) = 0$.

691 It remains to show that θ_∞ is not a saddle point for Lebesgue almost every initial value θ_0 . However,
 692 this is very similar to the proof in Lee et al. (2016). The only difference is that one would need to use
 693 a continuous-time version of the stable manifold theorem, which is readily available, for example in
 694 Bressan (2003). \square

695 Let $\{m_i\}_{i \in \mathbb{N}}$ denote the local minima of ℓ which are by assumption countable. Denote further by

$$\Theta_i := \{\theta_0 \in \mathbb{R}^J : \lim_{t \rightarrow \infty} \phi^t(\theta_0) \rightarrow m_i\} \quad (72)$$

696 the domain of attraction for the minimum m_i . The next theorem is then an easy consequence of
 697 Lemma 3.

698 **Theorem 5.** Assume that the loss function ℓ only has countably many local minima, is norm coe-
 699 rative, and satisfies the Lojasiewicz inequality. Let further $\theta_0 \sim Q_0$ for some $Q_0 \in \mathcal{P}(\mathbb{R}^J)$ such that
 700 $\sum_{i=1}^{\infty} Q_0(\Theta_i) = 1$. Then,

$$\phi^t(\theta_0) \xrightarrow{\mathcal{D}} \sum_{i=1}^{\infty} Q_0(\Theta_i) \delta_{m_i} =: Q_\infty \quad (73)$$

701 for $t \rightarrow \infty$. Here $\xrightarrow{\mathcal{D}}$ denotes convergence in distribution.

702 *Proof.* Let $\theta_0 \in \mathbb{R}^J$ be fixed. Due to Lemma 3, we know that

$$\phi^t(\theta_0) \rightarrow \sum_{i=1}^{\infty} m_i \mathbb{1}\{\theta_0 \in \Theta_i\} \quad (74)$$

703 for Lebesgue almost every θ_0 for $t \rightarrow \infty$. Here, $\mathbb{1}\{\cdot\}$ denotes the indicator function. Let Y now
 704 be a random variable with law Q_0 . By assumption, we know that $Y \in \Theta_i$ for some $i \in \mathbb{N}$ with
 705 probability 1. Hence,

$$\phi^t(Y) \rightarrow \sum_{i=1}^{\infty} m_i \mathbb{1}\{Y \in \Theta_i\} \quad (75)$$

706 almost surely for $t \rightarrow \infty$. Since almost sure convergence implies convergence in distribution, we
 707 conclude that

$$\phi^t(Y) \xrightarrow{\mathcal{D}} \mathcal{L}\left(\sum_{i=1}^{\infty} m_i \mathbb{1}\{Y \in \Theta_i\}\right), \quad (76)$$

708 where $\mathcal{L}(\cdot)$ denotes the law of a random variable. However, the law of the RHS is easily recognised
 709 as

$$\mathcal{L}\left(\sum_{i=1}^{\infty} m_i \mathbb{1}\{Y \in \Theta_i\}\right) = \sum_{i=1}^{\infty} Q_0(\Theta_i) \delta_{m_i}, \quad (77)$$

710 which concludes the proof. \square

711 **Remark 2.** Note that the condition

$$\sum_{i=1}^{\infty} Q_0(\Theta_i) = 1 \quad (78)$$

712 in Theorem 5 is easy to satisfy. According to Lemma 3 the set

$$\mathbb{R}^J \setminus \bigcup_{i=1}^n \Theta_i \quad (79)$$

713 has Lebesgue measure zero. Therefore, any Q_0 which has a density w.r.t. the Lebesgue measure
 714 will satisfy (78).

715 D Asymptotic distribution for deep Langevin ensembles

716 In this section, we analyse the objective

$$L(Q) := \int \ell(\theta) dQ(\theta) + \lambda \text{KL}(Q, P) \quad (80)$$

717 for $Q \in \mathcal{P}(\mathbb{R}^J)$. The corresponding particle method is given as:

- 718 • Sample $\theta_1(0), \dots, \theta_{N_E}(0)$ independently from Q_0 .
- 719 • Simulate the SDE $d\theta_n(t) = -\nabla V(\theta_n(t))dt + \sqrt{2\lambda}dB_n(t)$ for each $n = 1, \dots, N_E$.

720 Recall that $V(\theta) = \ell(\theta) - \lambda \log p(\theta)$. This case is well-studied in the literature and known as
 721 Langevin diffusion. Under mild assumptions (Chiang et al., 1987; Roberts and Tweedie, 1996),

$$\theta_n(t) \xrightarrow{\mathcal{D}} Q_{\infty} \quad (81)$$

722 for $t \rightarrow \infty$ and each particle $n = 1, \dots, N_E$ independently. The probability measure Q_{∞} has the
 723 density

$$q_{\infty}(\theta) = \frac{1}{Z} \exp\left(-\frac{V(\theta)}{\lambda}\right) \quad (82)$$

$$= \frac{1}{Z} \exp\left(-\frac{\ell(\theta)}{\lambda}\right) p(\theta), \quad (83)$$

724 where $Z > 0$ is the normalising constant. As a consequence, the WGF asymptotically produces
 725 samples from Q_{∞} . However, it is a priori unclear that Q_{∞} is in fact the same as the global minimiser
 726 Q^* of L .

727 We investigate this question by relating invariant measures to stationary points of the Wasserstein
 728 gradient.

729 **Definition 1.** (Liggett, 2010, Thm. 3.3.7) A measure Q is called an invariant measure (for a given
 730 Feller-process) if

$$\int Ah(\theta) dQ(\theta) = 0 \quad (84)$$

731 for all $h \in C_c^2(\mathbb{R}^J)$. Here A is the infinitesimal generator of the corresponding Feller-process.

732 Recall that the infinitesimal generator of the Langevin diffusion for $h \in C_c^2(\mathbb{R}^J)$ is given as

$$Ah = -\nabla V \cdot \nabla h + \lambda \Delta h. \quad (85)$$

733 **Definition 2.** A measure $Q \in \mathcal{P}_2(\mathbb{R}^J)$ is called a stationary point of the Wasserstein gradient if

$$\nabla_W L[Q](\theta) = 0 \quad (86)$$

734 for Q -almost every $\theta \in \mathbb{R}^J$.

735 In finite dimensions, it is well-known that a local minimiser is a stationary point of the gradient.
736 This carries over to the infinite-dimensional case, with a similar proof. Since we could not find this
737 result anywhere in the literature we included it for completeness.

738 **Lemma 4.** Let \widehat{Q} be a local minimiser of L , i.e. there exists and $\epsilon > 0$ such that

$$L(\widehat{Q}) \leq L(Q) \quad (87)$$

739 for all Q with $W_2(\widehat{Q}, Q) \leq \epsilon$. Then \widehat{Q} is a stationary point of the Wasserstein gradient in the sense
740 of Definition 2.

741 *Proof.* Let $h \in C_c^2(\mathbb{R}^J)$ be arbitrary and $\widehat{Q} \in \mathcal{P}_2(\mathbb{R}^J)$ be a local minimum of L . Further, let $\phi^t(\theta_0)$
742 be the solution to the initial value problem

$$\theta(0) = \theta_0 \quad (88)$$

$$\theta'(t) = \nabla h(\theta(t)) \quad (89)$$

743 for $t \in (-\epsilon, \epsilon)$ for some $\epsilon > 0$. We now define $Q(t) := \phi^t \# \widehat{Q}$ for $t \in (-\epsilon, \epsilon)$ where $f \# \mu$ denotes
744 the push-forward of the measures μ through the function f . In the Riemannian interpretation of the
745 Wasserstein space, $(Q(t))_{t \in (-\epsilon, \epsilon)}$ is a curve in $\mathcal{P}_2(\mathbb{R}^J)$ with tangent vector h at point \widehat{Q} (Ambrosio
746 et al., 2005, Chapter 8). We, further, define $f : (-\epsilon, \epsilon) \rightarrow \mathbb{R}$ as $f(t) := L(Q(t))$. Application of the
747 chain-rule (Ambrosio et al., 2005, p. 233) gives

$$f'(0) = \frac{d}{dt} L(Q(t)) \Big|_{t=0} \quad (90)$$

$$= \langle \nabla_W L[Q(0)], \nabla h \rangle_{L^2(Q(0))} \quad (91)$$

$$= \int \nabla_W L[\widehat{Q}](\theta) \cdot \nabla h(\theta) d\widehat{Q}(\theta). \quad (92)$$

748 We know that f has a local minimum at $t = 0$ and, therefore, $f'(0) = 0$ which gives

$$0 = \int \nabla_W L[\widehat{Q}](\theta) \cdot \nabla h(\theta) d\widehat{Q}(\theta). \quad (93)$$

749 Since (93) holds for arbitrary test functions $h \in C_c^2(\mathbb{R}^J)$ and as $C_c^2(\mathbb{R}^J)$ is dense in $L^2(\widehat{Q})$, we
750 obtain that $\nabla_W L[\widehat{Q}](\theta) = 0$ for \widehat{Q} -a.e $\theta \in \mathbb{R}^J$. \square

751 The next lemma relates invariant measures and stationary points of the Wasserstein gradient for
752 infinitesimal generators of the form (85). It will prove extremely useful to translate between the
753 Langevin diffusion literature and our optimisation perspective.

754 **Lemma 5.** Let $Q \in \mathcal{P}_2(\mathbb{R}^J)$ be such that Q has a density q with respect to the Lebesgue measure.
755 Then, the following two statements are equivalent:

- 756 • Q is a stationary point of the Wasserstein gradient.
- 757 • Q is an invariant measure.

758 *Proof.* Let Q be a measure with density q . Recall that the generator of the Langevin diffusion is for
759 $h \in C_c^2(\mathbb{R}^J)$ given as

$$Ah = -\nabla V \cdot \nabla h + \lambda \Delta h. \quad (94)$$

760 By partial integration, it is easy to verify that the L^2 - adjoint (w.r.t the Lebesgue measure) is given
761 as

$$A^*h = \nabla \cdot (h \cdot \nabla V) + \lambda \Delta h. \quad (95)$$

762 We, therefore, conclude that

$$\int Ah(\theta) dQ(\theta) = \int Ah(\theta)q(\theta) d\theta \quad (96)$$

$$= \int h(\theta)A^*q(\theta) d\theta \quad (97)$$

$$= \int h(\theta) \left(\nabla \cdot (q(\theta) \cdot \nabla V(\theta)) + \lambda \Delta q(\theta) \right) d\theta. \quad (98)$$

763 Furthermore, we have $\nabla_W L[Q] = \nabla V + \lambda \nabla \log q$, and therefore

$$\int \nabla_W L[Q](\theta) \cdot \nabla h(\theta) dQ(\theta) = \int \nabla_W L[Q](\theta) \cdot \nabla h(\theta)q(\theta) d\theta \quad (99)$$

$$= \int \left(\nabla V(\theta)q(\theta) + \lambda \nabla q(\theta) \right) \cdot \nabla h(\theta) d\theta \quad (100)$$

$$= - \int h(\theta) \left(\nabla \cdot (q(\theta) \nabla V(\theta)) + \lambda \Delta q(\theta) \right) d\theta, \quad (101)$$

764 where the last line follows from applying partial integration. This allows us to conclude that

$$\int Ah(\theta) dQ(\theta) = - \int \nabla_W L[Q](\theta) \cdot \nabla h(\theta) dQ(\theta) \quad (102)$$

765 whenever Q has a density. As a consequence we have that Q is invariant if and only if it is a
766 stationary point of the Wasserstein gradient. \square

767 Lemma 5 allows us to move between the optimisation and stochastic differential equation perspec-
768 tive. In Appendix A we discussed the existence and uniqueness of a global minimiser Q^* of L .
769 We know that Q^* has a density since the Kullback-Leibler divergence would be infinite otherwise
770 (assuming P has a Lebesgue-density which we assume throughout the paper). Lemma 4 guarantees
771 that Q^* is a stationary point of the Wasserstein gradient. Due to Lemma 5, we can infer that Q^*
772 must be an invariant measure. However, due to the uniqueness of the invariant measure under the
773 previously mentioned mild assumptions (Chiang et al., 1987; Roberts and Tweedie, 1996), we can
774 conclude that $Q^* = Q_\infty$.

775 E Asymptotic distribution of deep repulsive Langevin ensembles

776 In this section, we consider

$$L(Q) = \int \ell(\theta) dQ(\theta) + \frac{\lambda_1}{2} \text{MMD}(Q, P)^2 + \lambda_2 \text{KL}(Q, P) \quad (103)$$

$$= \int V(\theta) dQ(\theta) + \frac{\lambda_1}{2} \int \kappa(\theta, \theta') dQ(\theta)dQ(\theta') - \lambda_2 H(Q) + \text{const}, \quad (104)$$

777 as optimisation objective. Here, $H(Q) = - \int \log q(\theta)q(\theta) d\theta$ denotes the differential entropy.

778 Recall that in this case $V(\theta) = \ell(\theta) - \lambda_1 \mu_P(\theta) - \lambda_2 \log p(\theta)$. We already discussed in Appendix B
779 that the McKean-Vlasov process of the form

$$\theta(0) \sim Q_0 \quad (105)$$

$$d\theta(t) = - \left(\nabla V(\theta(t)) + \lambda_1 (\nabla_1 \kappa * Q_t)(\theta(t)) \right) dt + \sqrt{2\lambda_2} dB(t), \quad (106)$$

780 with $(B(t))_{t \geq 0}$ being a Brownian motion achieves the desired density evolution. Furthermore, the
781 particle approximation of (105) is given as

$$d\theta_n(t) = - \left(\nabla V(\theta_n(t)) + \frac{\lambda_1}{N_E} \sum_{j=1}^{N_E} (\nabla_1 \kappa)(\theta_n(t), \theta_j(t)) \right) dt + \sqrt{2\lambda_2} dB_n(t) \quad (107)$$

782 for $n = 1, \dots, N_E$ where $N_E \in \mathbb{N}$ denotes the number of particles.

783 The approach follows the same procedure as in Appendix [D](#). We show the notions of invariant
784 measures and stationary points of the Wasserstein gradient are the same for measures with Lebesgue
785 density. We start by introducing the concept of an invariant measure for a nonlinear Markov process
786 ([Ahmed and Ding, 1993](#), Definition 1).

787 **Definition 3.** A measure Q is called an invariant measure for a nonlinear Markov process with the
788 family of infinitesimal generators $\{A[Q]\}_{Q \in \mathcal{P}(\mathbb{R}^J)}$ if

$$\int A[Q]h(\theta) dQ(\theta) = 0 \quad (108)$$

789 for all $h \in C_c^2(\mathbb{R}^J)$.

790 Recall that the family of infinitesimal generators in our case is given as

$$(A[Q]h)(\theta) := -\left(\nabla V(\theta) + \lambda_1(\nabla_1 \kappa * Q)(\theta)\right) \cdot \nabla h(\theta) + \lambda_2 \Delta h. \quad (109)$$

791 for $h \in C_c^2(\mathbb{R}^J, \mathbb{R})$. In analogy to Lemma [5](#), we obtain the following result.

792 **Lemma 6.** Let $Q \in \mathcal{P}_2(\mathbb{R}^J)$ be such that Q has a density q with respect to the Lebesgue measure.
793 Then, the following two statements are equivalent:

- 794 • Q is a stationary point of the Wasserstein gradient for L in [\(103\)](#) in the sense of Def. [2](#)
- 795 • Q is an invariant measure for the McKean-Vlasov process with infinitesimal generator
796 defined in [\(109\)](#)

797 *Proof.* First, we notice that

$$\int A[Q]h(\theta) dQ(\theta) = \int A[Q]h(\theta)q(\theta) d\theta \quad (110)$$

$$= \int h(\theta) (A^*[Q]q)(\theta) d\theta. \quad (111)$$

798 Recall, that $A^*[Q]$ denotes the L^2 -adjoint of the operator $A[Q]$ and that it is given as

$$(A^*[Q]h)(\theta) = \nabla \cdot \left(h(\theta) (\nabla V(\theta) + \lambda_1(\nabla_1 \kappa * Q)(\theta) + \lambda_2 \nabla \log(h(\theta))) \right) \quad (112)$$

799 for $h \in C^2(\mathbb{R}^J, \mathbb{R})$ with compact support. This implies

$$(A^*[Q]q)(\theta) = \nabla \cdot \left(q(\theta) (\nabla V(\theta) + \lambda_1(\nabla_1 \kappa * Q)(\theta)) \right) + \lambda_2 \Delta q(\theta). \quad (113)$$

800 We plug this into [\(111\)](#) to obtain

$$\int A[Q]h(\theta) dQ(\theta) \quad (114)$$

$$= \int h(\theta) \nabla \cdot \left(q(\theta) (\nabla V(\theta) + \lambda_1(\nabla_1 \kappa * Q)(\theta)) \right) d\theta + \int \lambda_2 h(\theta) \Delta q(\theta) d\theta. \quad (115)$$

801 On the other hand, we have that

$$\nabla_W L[Q](\theta) = \nabla V(\theta) + \lambda_1(\nabla_1 \kappa * Q)(\theta) + \lambda_2 \nabla \log q(\theta), \quad (116)$$

802 and therefore

$$\int \nabla L[Q](\theta) \cdot \nabla h(\theta) dQ(\theta) \quad (117)$$

$$= \int \left(\nabla V(\theta) + \lambda_1(\nabla_1 \kappa * Q)(\theta) + \lambda_2 \nabla \log q(\theta) \right) \cdot \nabla h(\theta) dQ(\theta) \quad (118)$$

$$= \int q(\theta) (\nabla V(\theta) + \lambda_1(\nabla_1 \kappa * Q)(\theta)) \cdot \nabla h(\theta) d\theta + \lambda_2 \int \nabla q(\theta) \cdot \nabla h(\theta) d\theta \quad (119)$$

$$= - \int \nabla \cdot \left(q(\theta) (\nabla V(\theta) + \lambda_1(\nabla_1 \kappa * Q)(\theta)) \right) h(\theta) d\theta - \lambda_2 \int q(\theta) \Delta h(\theta) d\theta, \quad (120)$$

803 where the last line follows from partial integration. Comparing (115) to (120) gives

$$\int A[Q]h(\theta) dQ(\theta) = - \int \nabla L[Q](\theta) \cdot \nabla h(\theta) dQ(\theta) \quad (121)$$

804 for all $h \in C_c^2(\mathbb{R}^J)$ whenever Q has a density. This immediately implies that Q is invariant iff it is
805 a stationary point. \square

806 Again, we leverage this correspondence between stationary point and invariant measures. There is
807 a rich literature on ergodicity of nonlinear Markov processes. For example, Theorem 2 of Vereten-
808 nikov (2006) specifies conditions on κ and V such that

$$Q^{n, N_E}(t) \xrightarrow{\mathcal{D}} Q_\infty \quad (122)$$

809 for $N_E, t \rightarrow \infty$. Here $Q^{n, N_E}(t)$ denotes the law of a fixed particle $\theta_n(t)$, $n = 1, \dots, N_E$, whose
810 distribution is characterised by the SDE (59). The measure Q_∞ is the unique invariant measure of
811 the nonlinear Markov process. By Lemma 6 every invariant measure is a stationary point of the
812 Wasserstein gradient and vice versa. Hence, existence and uniqueness of the stationary point of the
813 Wasserstein gradient is immediately implied. However, since the global minimiser Q^* is a stationary
814 point of the Wasserstein gradient (cf. Lemma 4), we conclude by uniqueness that $Q_\infty = Q^*$.

815 F Asymptotic analysis of deep repulsive ensembles

816 In this section, we consider the objective

$$L(Q) := \int \ell(\theta) dQ(\theta) + \lambda \text{MMD}(Q, P) \quad (123)$$

817 for $Q \in \mathcal{P}(\mathbb{R}^J)$. The corresponding McKean-Vlasov process is of the form

$$d\theta(t) = - \left(\nabla V(\theta(t)) + \lambda (\nabla_1 \kappa * Q_t)(\theta(t)) \right) dt, \quad (124)$$

818 where Q_t denotes the distribution of $\theta(t)$ and $V(\theta) = \ell(\theta) - \mu_P(\theta)$ with $\mu_P(\theta) = \int \kappa(\theta, \theta') dP(\theta')$
819 the kernel mean-embedding of P . We call the particle method in this case **deep repulsive ensembles**
820 **(DRE)**.

821 The existence of the global minimiser Q^* is still guaranteed under the assumptions in Appendix A.
822 Lemma 4 guarantees that Q^* is a stationary point of the Wasserstein gradient, i.e.

$$\nabla V(\theta) + \lambda (\nabla_1 \kappa * Q^*)(\theta) = 0 \quad (125)$$

823 for Q^* -a.e. $\theta \in \mathbb{R}^J$. Recall that the infinitesimal generator in this case is given as

$$(A[Q]h)(\theta) := - \left(\nabla V(\theta) + \lambda (\nabla_1 \kappa * Q)(\theta) \right) \cdot \nabla h(\theta) \quad (126)$$

824 for $Q \in \mathcal{P}(\mathbb{R}^J)$, $h \in C_c^2(\mathbb{R}^J)$. It immediately follows from the definition that

$$(A[Q]h)(\theta) = - \nabla_W L[Q](\theta) \cdot \nabla h(\theta) \quad (127)$$

825 for all $h \in C_c^2(\mathbb{R}^J)$, $\theta \in \mathbb{R}^J$. As in Lemma 5 & 6, this implies that each stationary point of
826 the Wasserstein gradient is an invariant measure of the McKean-Vlasov process and vice versa. In
827 Appendix D & E, we cite relevant literature that guarantees uniqueness of the invariant measure,
828 which is a necessary (but not sufficient) condition for convergence to the invariant measure. The
829 next theorem shows that uniqueness will in general not hold without the presence of the diffusion
830 term.

831 **Theorem 6.** *The invariant measure for the McKean-Vlasov process with the family of generators*
832 *$(A[Q])_{Q \in \mathcal{P}(\mathbb{R}^J)}$ defined in (127) is (in general) not unique.*

833 *Proof.* Let $N_E \in \mathbb{N}$ and define $\tilde{L} : (\mathbb{R}^J)^{N_E} \rightarrow \mathbb{R}$ as

$$\tilde{L}(\theta_1, \dots, \theta_{N_E}) := \sum_{i=1}^{N_E} V(\theta_i) + \frac{\lambda}{2N_E} \sum_{i,j=1}^{N_E} \kappa(\theta_i, \theta_j). \quad (128)$$

834 Assume that V is bounded from below and norm-coercive. Then \tilde{L} is bounded from below and
 835 norm-coercive and therefore we can find a global minimiser $\theta^* := (\theta_1^*, \dots, \theta_{N_E}^*) \in (\mathbb{R}^J)^{N_E}$ of \tilde{L} .
 836 Since \tilde{L} is differentiable, we know that θ^* is a stationary point of the gradient which implies

$$\nabla V(\theta_i^*) + \frac{\lambda}{N_E} \sum_{j=1}^{N_E} (\nabla_1 \kappa)(\theta_i^*, \theta_j^*) = 0 \quad (129)$$

837 for all $i = 1, \dots, N_E$. Here, we assume that the kernel κ is symmetric, which is standard in the
 838 MMD literature. Note that (129) is equivalent to

$$\nabla V(\theta) + \lambda(\nabla_1 \kappa * \hat{Q})(\theta) = 0 \quad (130)$$

839 for \hat{Q} -a.e. $\theta \in \mathbb{R}^J$ where

$$\hat{Q}(d\theta) := \frac{1}{N_E} \sum_{j=1}^{N_E} \delta_{\theta_j^*}(d\theta). \quad (131)$$

840 This means that \hat{Q} is a stationary point of the Wasserstein gradient, and therefore an invariant mea-
 841 sure for the McKean-Vlasov process. Since $N_E \in \mathbb{N}$ was arbitrary, we have constructed countably
 842 many invariant measures and therefore uniqueness can't hold in general. \square

843 The reason that non-uniqueness of the invariant measure is an immediate contradiction to conver-
 844 gence is the following: If we initialise with any of the invariant measures constructed in the proof of
 845 Theorem 6, then the particle distribution of the McKean-Vlasov process will remain unchanged over
 846 time. Convergence to the global minimiser can therefore surely not hold for arbitrary initialisation
 847 Q_0 . It may be possible to construct conditions on Q_0 under which convergence still holds. For
 848 example, for Stein variational gradient descent a similar issue occurs. However, in this case one can
 849 guarantee convergence (Lu et al., 2019, Theorem 2.8) if Q_0 has a Lebesgue-density (and if the ker-
 850 nel satisfies further restrictive assumptions). The existence of conditions that guarantee convergence
 851 for DRE remains an open problem.

852 G Implementation details

853 In Appendix A we derived the following algorithm:

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

855 **Step 2:** Evolve the particles forward in time according to

$$\theta_{n,k+1} = \theta_{n,k} - \eta \left(\nabla V(\theta_{n,k}) + \frac{\lambda_1}{N_E} \sum_{j=1}^{N_E} (\nabla_1 \kappa)(\theta_{n,k}, \theta_{j,k}) \right) + \sqrt{2\eta\lambda_2} Z_{n,k} \quad (132)$$

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

857 We can generate samples from DE, DLE and DRLE by setting the potential and regularisation
 858 parameters as described below:

- 859 • Deep ensembles: $V(\theta) = \ell(\theta)$, $\lambda_1 = 0$, $\lambda_2 = 0$
- 860 • Deep Langevin ensembles: $V(\theta) = \ell(\theta) - \lambda \log p(\theta)$, $\lambda_1 = 0$, $\lambda = \lambda_2$
- 861 • Deep repulsive Langevin ensembles: $V(\theta) = \ell(\theta) - \lambda_1 \log p(\theta) - \lambda_2 \mu_P(\theta)$

862 Due to Appendix D & E, we can think of $\theta_{1,K}, \dots, \theta_{N_E,K}$ as approximately sampled from the
 863 global minimiser Q^* for DLE and DRLE if K is large enough. All experiments use the SE kernel
 864 given as

$$\kappa(\theta, \theta') = \exp \left(- \frac{\|\theta - \theta'\|^2}{2\sigma_\kappa^2} \right) \quad (133)$$

865 with lengthscale parameter $\sigma_\kappa > 0$. The kernel mean embedding μ_P can easily be approximated as

$$\mu_P(\theta) = \frac{1}{M} \sum_{i=1}^M \kappa(\theta, \theta_i), \quad \theta \in \mathbb{R}^J, \quad (134)$$

866 where $\theta_1, \dots, \theta_M \sim P$ independently. We chose $M = 20$.

867 **G.1 Toy example: global minimiser**

868 We describe details regarding the experiments conducted to produce Figure 2 below.

869 We generate $N_E = 300$ particles and make the following choices:

- 870 • Loss: $\ell(\theta) := \frac{3}{2}(\frac{1}{4}\theta^4 + \frac{1}{3}\theta^3 - \theta^2) - \frac{3}{8}$
- 871 • Prior: $P \sim \mathcal{N}(0, 1)$ and therefore $\log p(\theta) = -\frac{1}{2}\theta^2$
- 872 • Initialisation: $Q_0 = P$
- 873 • Reg. parameter: $\lambda_{DLE} = 1, \lambda_{DRLE} = 1, \lambda'_{DRLE} = 1$
- 874 • Step size: $\eta = 10^{-4}$, Iterations: $K = 100,000$
- 875 • Kernel lengthscale, σ_κ , is chosen according to the median heuristic (Garreau et al., 2017)
- 876 based on samples from the prior P

877 The loss is constructed such that we have a global minimum at $\theta = -2$, a turning point at $\theta = 0$,
878 and a local minimum at $\theta = 1$.

- 879 • Deep ensembles: The optimal Q^* is a Dirac measure located at the global minimiser $\theta =$
880 -2 . However, as we proved in Theorem 1 the WGF produce samples from

$$Q_\infty(d\theta) = \frac{1}{2}\delta_{-2}(d\theta) + \frac{1}{2}\delta_1(d\theta), \quad (135)$$

881 as $(-\infty, 0)$ is the region of attraction for the global minimum and $(0, \infty)$ for the local
882 minimum which both have probability 0.5 under $Q_0 = P = \mathcal{N}(0, 1)$. In particular, $Q_\infty \neq$
883 Q^* as expected.

- 884 • Deep Langevin ensembles: The optimal measure has the pdf

$$q^*(\theta) \propto \exp\left(-\frac{\ell(\theta)}{\lambda}\right)p(\theta) \quad (136)$$

885 for $\theta \in \mathbb{R}$. As expected the WGF produces samples from Q^* .

- 886 • Deep repulsive Langevin ensembles: The optimal q^* for deep repulsive ensembles is harder
887 to determine. From the condition that q^* is a stationary point of the Wasserstein gradient,
888 we can derive that $u(\theta) := \log q^*(\theta)$ satisfies the integro-differential equation

$$u'(\theta) = -\frac{1}{\lambda_2}V'(\theta) - \frac{\lambda_1}{\lambda_2} \int (\nabla_1 \kappa)(\theta, \theta') \exp(u(\theta')) d\theta' \quad (137)$$

889 with some initial value $u(0) = u_0$. In principle, we could choose u_0 such that $q(\theta) :=$
890 $\exp(u(\theta))$ integrates to 1. However, since we do not know the appropriate initial condition
891 a priori, we choose an arbitrary u_0 and normalise the pdf afterwards. We use an numerical
892 solver to evaluate $u(\theta)$ on a fixed grid. As expected, the WGF produces samples from Q^*
893 in this case.

894 **G.2 Toy example: multimodal loss**

895 The details below correspond to the experimental results presented in Figure 3

896 **DE, DLE, DRLE** We generate $N_E = 300$ particles and make the following choices:

- 897 • Loss: $\ell(\theta) = -\log \sum_{i=1}^4 \frac{1}{4}\mathcal{N}(\theta; \mu_i, I_2)$, $\theta \in \mathbb{R}^2$, $\mu_i = (\pm 3, \pm 3)^T$, $i = 1, \dots, 4$
- 898 • Prior: P flat and therefore $\log p(\theta) = 0$
- 899 • Initialisation: $Q_0 \sim \mathcal{N}(0, I_2)$
- 900 • Reg. parameter: $\lambda_{DLE} = 0.2, \lambda_{DRLE} = 0.2, \lambda'_{DRLE} = 0.6$
- 901 • Step size: $\eta = 0.1$, Iterations: $K = 10,000$
- 902 • Kernel lengthscale, σ_κ , is chosen according to the median heuristic (Garreau et al., 2017)
- 903 based on samples from the prior P

904 Note that for a translation-invariant kernel such as the SE kernel we obtain for the flat prior P that

$$\mu_P(\theta) = \int_{-\infty}^{\infty} \kappa(\theta, \theta') d\theta' \quad (138)$$

$$= \int_{-\infty}^{\infty} \phi(\theta - \theta') d\theta' \quad (139)$$

$$= \int_{-\infty}^{\infty} \phi(\xi) d\xi, \quad (140)$$

905 where the second line follows from the fact that we can write any translation-invariant kernel as
 906 $\kappa(\theta, \theta') = \phi(\theta - \theta')$ for some function $\phi : \mathbb{R}^J \rightarrow \mathbb{R}$ and the second line is simple variable substitution.
 907 If (140) is finite, the above expression is well-defined and therefore μ_P constant. Note that in
 908 particular for the SE kernel, we have $\phi(\xi) = \exp(-\|\xi\|^2/(2\sigma_\kappa^2))$ and therefore (140) is finite. As a
 909 consequence, we have that for a flat prior P the gradient of the potential V is the same for all three
 910 methods. This means that the loss ℓ isn't adjusted and the only difference between the three methods
 911 is the presence of repulsion and noise effects.

912 **Remark 3.** The astute reader may have noticed that a flat prior P is in fact not covered by our theory
 913 in Appendix A. The problem is that $\text{KL}(\cdot, \mathcal{L})$, where \mathcal{L} denotes the Lebesgue measure, is not positive
 914 (and not even bounded from below). To see this, choose $Q = \mathcal{N}(0, \Sigma)$ with $\Sigma = \text{diag}(\sigma_1^2, \sigma_2^2)$ and
 915 note that

$$\text{KL}(Q, \mathcal{L}) = \int \log q(\theta) q(\theta) d\theta = -\text{H}(Q), \quad (141)$$

916 where $\text{H}(Q)$ denotes the differential entropy. For a Gaussian, it is known that

$$\text{H}(Q) = \frac{1}{2} \log((2\pi e)^J \det(\Sigma)) = \frac{1}{2} (\log(2\pi e)^J + \log(\sigma_1^2) + \log(\sigma_2^2)) \quad (142)$$

917 and therefore if either $\sigma_1^2 \rightarrow \infty$ or $\sigma_2^2 \rightarrow \infty$ then $\text{KL}(Q, \mathcal{L}) \rightarrow -\infty$. However, note that this
 918 difficulty is rather technical in nature and can easily be remedied. Instead of \mathcal{L} , we could have chosen
 919 the uniform prior $P \sim U(-10^{100}, 10^{100})$. In this case, the positivity of $\text{KL}(\cdot, P)$ is guaranteed by
 920 Jensen's inequality. This choice of P gives—up to an additive constant—the same objective as
 921 a flat prior and up to machine precision the same kernel mean embedding μ_P . It is, therefore,
 922 algorithmically irrelevant if P is flat or uniform on a very large set.

923 **FD-GVI** We use the same prior and loss as for DE, DLE and DRLE. We parameterise the variational
 924 family as independent Gaussian, i.e.

$$\mathcal{Q} = \{\mathcal{N}(\mu, \Sigma) \mid \mu \in \mathbb{R}^2, \Sigma = \text{diag}(\exp(\beta_1), \exp(\beta_2)), \beta := (\beta_1, \beta_2)^2 \in \mathbb{R}^2\}. \quad (143)$$

925 We learn the variational parameters $\nu := (\mu, \beta) \in \mathbb{R}^4$ by minimising

$$\tilde{L}(\nu) = \int \ell(\theta) dQ_\nu(\theta) + \lambda \text{KL}(Q_\nu, P) \quad (144)$$

$$= \int \ell(\theta) dQ_\nu(\theta) - \lambda H(\mathcal{N}(\mu, \Sigma)) \quad (145)$$

$$\approx \frac{1}{200} \sum_{j=1}^{200} \ell(\mu + \Sigma^{0.5} Z_j) - \frac{\lambda}{2} \log((2\pi e)^2 \exp(\beta_1) \exp(\beta_2)) \quad (146)$$

$$= \frac{1}{200} \sum_{j=1}^{200} \ell(\mu + \Sigma^{0.5} Z_j) - \frac{\lambda}{2} (\beta_1 + \beta_2) + \text{const}, \quad (147)$$

926 where $Z_1, \dots, Z_{200} \sim \mathcal{N}(0, I_2)$ and $H(\mathcal{N}(\mu, \Sigma))$ denotes the differential entropy of the normal
 927 distribution. For the regularisation parameter, we chose $\lambda = 0.5$.

928 G.3 Toy example: more modes than particles

929 We generate $N_E = 20$ particles and make the following choices:

- 930 • Loss: $\ell(\theta) = -|\sin(\theta)|$, $\theta \in [-M\pi, M\pi]$, with $M = 1000$
- 931 • Prior: P flat and therefore $\log p(\theta) = 0$ and $\mu_P = \text{const.}$ (cf. Appendix [G.2](#))
- 932 • Initialisation: $Q_0 \sim U(-M\pi, M\pi)$
- 933 • Reg. parameter: $\lambda_{DLE} = 0.001$, $\lambda_{DRLE} = 0.001$, $\lambda'_{DRLE} = 0.6$
- 934 • Step size: $\eta = 0.01$, Iterations: $K = 1.000$
- 935 • Kernel lengthscale, σ_κ , is chosen according to the median heuristic ([Garreau et al., 2017](#))
- 936 based on samples from the prior P

937 Note that ℓ has $2M = 2000$ local minima at locations

$$m_i := \frac{\pi}{2} + i\pi, \quad i \in \{-M, \dots, 0, \dots, (M-1)\}. \quad (148)$$

938 Due to the flat prior $\nabla V = \nabla \ell$ for all three methods. We observe that it is hard to distinguish the
939 methods since most particles are in their local modes by themselves.

940 G.4 UCI Regression

941 The UCI data sets are licensed under Creative Commons Attribution 4.0 International license (CC
942 BY 4.0). Following [Lakshminarayanan et al. \(2017\)](#), we train 5 one-hidden-layer neural networks
943 f_θ with 50 hidden nodes for 40 epochs. We split each data set into train (81% of samples), validation
944 (9% of samples), and test set (10% of samples). Based on the best hyperparameter runs (according
945 to a Gaussian NLL) found via grid search on a validation data set, we make the following choices:

- 946 • Loss: $\ell(\theta) = \frac{1}{N} \sum_{n=1}^N (f_\theta(x_n) - y_n)^2$ where $\{x_n, y_n\}_{n=1}^N$ are paired observations.
- 947 • Prior: $P \sim \mathcal{N}(0, 1)$
- 948 • Initialisation: Kaiming intilisation, i.e. for each layer $l \in \{1, \dots, L\}$ that maps features with
949 dimensionality n_{l-1} into dimensionality n_l , we sample $Q_{l,0} \sim \mathcal{N}(0, 2/n_l)$
- 950 • Reg. parameter: $\lambda_{DLE} = 10^{-4}$, $\lambda_{DRLE} = 10^{-4}$, $\lambda'_{DRLE} = 10^{-2}$
- 951 • Step size: $\eta = 0.1$, Iterations: $K = 10,000$
- 952 • Kernel lengthscale, σ_κ , is chosen according to the median heuristic ([Garreau et al., 2017](#))
- 953 based on samples from the prior P

954 G.5 Compute

955 While the final experimental results can be run within approximately an hour on a single GeForce
956 RTX 3090 GPU, the complete compute needed for the final results, debugging runs, and sweeps
957 amounts to around 9 days.

958 References

- 959 Ahmed, N. and Ding, X. (1993). On invariant measures of nonlinear Markov processes. *Journal of*
960 *Applied Mathematics and Stochastic Analysis*, 6(4):385–406.
- 961 Alquier, P. (2021a). Non-exponentially weighted aggregation: regret bounds for unbounded loss
962 functions. In *International Conference on Machine Learning*, pages 207–218. PMLR.
- 963 Alquier, P. (2021b). User-friendly introduction to PAC-Bayes bounds. *arXiv preprint*
964 *arXiv:2110.11216*.
- 965 Alquier, P. and Guedj, B. (2018). Simpler PAC-Bayesian bounds for hostile data. *Machine Learning*,
966 107(5):887–902.
- 967 Altamirano, M., Briol, F-X., and Knoblauch, J. (2023). Robust and scalable Bayesian online
968 changepoint detection. *arXiv preprint arXiv:2302.04759*.
- 969 Ambrosio, L., Gigli, N., and Savaré, G. (2005). *Gradient flows: in metric spaces and in the space*
970 *of probability measures*. Springer Science & Business Media.

- 971 Arbel, M., Korba, A., Salim, A., and Gretton, A. (2019). Maximum mean discrepancy gradient flow.
972 *Advances in Neural Information Processing Systems*, 32.
- 973 Barbu, V. and Röckner, M. (2020). From nonlinear Fokker–Planck equations to solutions of distri-
974 bution dependent sde. *arXiv preprint arXiv:1808.10706*.
- 975 Bégin, L., Germain, P., Laviolette, F., and Roy, J.-F. (2016). PAC-Bayesian bounds based on the
976 Rényi divergence. In *Artificial Intelligence and Statistics*, pages 435–444.
- 977 Bissiri, P. G., Holmes, C. C., and Walker, S. G. (2016). A general framework for updating be-
978 lief distributions. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*,
979 78(5):1103–1130.
- 980 Blundell, C., Cornebise, J., Kavukcuoglu, K., and Wierstra, D. (2015). Weight uncertainty in neural
981 network. In *International conference on machine learning*, pages 1613–1622. PMLR.
- 982 Bressan, A. (2003). Tutorial on the center manifold theorem. *Hyperbolic systems of balance laws*,
983 1911:327–344.
- 984 Chiang, T.-S., Hwang, C.-R., and Sheu, S. J. (1987). Diffusion for global optimization in R^n . *SIAM*
985 *Journal on Control and Optimization*, 25(3):737–753.
- 986 Colding, T. H. and Minicozzi II, W. P. (2014). Lojasiewicz inequalities and applications. *arXiv*
987 *preprint arXiv:1402.5087*.
- 988 D’Angelo, F. and Fortuin, V. (2021). Repulsive deep ensembles are Bayesian. *Advances in Neural*
989 *Information Processing Systems*, 34:3451–3465.
- 990 Ermak, D. L. (1975). A computer simulation of charged particles in solution. i. technique and
991 equilibrium properties. *The Journal of Chemical Physics*, 62(10):4189–4196.
- 992 Fort, S., Hu, H., and Lakshminarayanan, B. (2019). Deep ensembles: A loss landscape perspective.
993 *arXiv preprint arXiv:1912.02757*.
- 994 Gal, Y. and Ghahramani, Z. (2016). Dropout as a Bayesian approximation: Representing model
995 uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059.
996 PMLR.
- 997 Garreau, D., Jitkrittum, W., and Kanagawa, M. (2017). Large sample analysis of the median heuristic.
998 *arXiv preprint arXiv:1707.07269*.
- 999 Glaser, P., Arbel, M., and Gretton, A. (2021). Kale flow: A relaxed kl gradient flow for probabilities
1000 with disjoint support. *Advances in Neural Information Processing Systems*, 34:8018–8031.
- 1001 Graves, A. (2011). Practical variational inference for neural networks. *Advances in neural informa-*
1002 *tion processing systems*, 24.
- 1003 Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B., and Smola, A. (2012). A kernel two-
1004 sample test. *Journal of Machine Learning Research*, 13(25):723–773.
- 1005 Grünwald, P. (2011). Safe learning: bridging the gap between Bayes, MDL and statistical learn-
1006 ing theory via empirical convexity. In *Proceedings of the 24th Annual Conference on Learning*
1007 *Theory*, pages 397–420.
- 1008 Guedj, B. and Shawe-Taylor, J. (2019). A primer on pac-Bayesian learning. In *ICML 2019-Thirty-*
1009 *sixth International Conference on Machine Learning*.
- 1010 Haddouche, M. and Guedj, B. (2023). Wasserstein PAC-Bayes learning: A bridge between general-
1011 isation and optimisation. *arXiv preprint arXiv:2304.07048*.
- 1012 Husain, H. and Knoblauch, J. (2022). Adversarial interpretation of Bayesian inference. In *Interna-*
1013 *tional Conference on Algorithmic Learning Theory*, pages 553–572. PMLR.
- 1014 Izmailov, P., Vikram, S., Hoffman, M. D., and Wilson, A. G. G. (2021). What are Bayesian neural
1015 network posteriors really like? In *International conference on machine learning*, pages 4629–
1016 4640. PMLR.

- 1017 Jewson, J., Smith, J., and Holmes, C. (2018). Principles of Bayesian inference using general diver-
1018 gence criteria. *Entropy*, 20(6):442.
- 1019 Knoblauch, J. (2019). Frequentist consistency of generalized variational inference. *arXiv preprint*
1020 *arXiv:1912.04946*.
- 1021 Knoblauch, J. (2021). *Optimization-centric generalizations of Bayesian inference*. PhD thesis,
1022 University of Warwick.
- 1023 Knoblauch, J., Jewson, J., and Damoulas, T. (2018). Doubly robust Bayesian inference for non-
1024 stationary streaming data using β -divergences. In *Advances in Neural Information Processing*
1025 *Systems (NeurIPS)*, pages 64–75.
- 1026 Knoblauch, J., Jewson, J., and Damoulas, T. (2022). An optimization-centric view on Bayes’
1027 rule: Reviewing and generalizing variational inference. *Journal of Machine Learning Research*,
1028 23(132):1–109.
- 1029 Kolokoltsov, V. N. (2010). *Nonlinear Markov processes and kinetic equations*, volume 182. Cam-
1030 bridge University Press.
- 1031 Korba, A., Aubin-Frankowski, P.-C., Majewski, S., and Ablin, P. (2021). Kernel Stein discrepancy
1032 descent. In *International Conference on Machine Learning*, pages 5719–5730. PMLR.
- 1033 Lakshminarayanan, B., Pritzel, A., and Blundell, C. (2017). Simple and scalable predictive un-
1034 certainty estimation using deep ensembles. *Advances in neural information processing systems*,
1035 30.
- 1036 Lee, J. D., Simchowitz, M., Jordan, M. I., and Recht, B. (2016). Gradient descent only converges to
1037 minimizers. In *Conference on learning theory*, pages 1246–1257. PMLR.
- 1038 Lichman, M. (2013). UCI machine learning repository.
- 1039 Liggett, T. M. (2010). *Continuous time Markov processes: an introduction*, volume 113. American
1040 Mathematical Soc.
- 1041 Liu, Q. (2017). Stein variational gradient descent as gradient flow. *Advances in neural information*
1042 *processing systems*, 30.
- 1043 Liu, Q. and Wang, D. (2016). Stein variational gradient descent: A general purpose Bayesian
1044 inference algorithm. *Advances in neural information processing systems*, 29.
- 1045 Louizos, C. and Welling, M. (2017). Multiplicative normalizing flows for variational Bayesian
1046 neural networks. In *International Conference on Machine Learning*, pages 2218–2227. PMLR.
- 1047 Lu, J., Lu, Y., and Nolen, J. (2019). Scaling limit of the stein variational gradient descent: The mean
1048 field regime. *SIAM Journal on Mathematical Analysis*, 51(2):648–671.
- 1049 Matsubara, T., Knoblauch, J., Briol, F.-X., and Oates, C. J. (2022). Robust generalised Bayesian
1050 inference for intractable likelihoods. *Journal of the Royal Statistical Society Series B: Statistical*
1051 *Methodology*, 84(3):997–1022.
- 1052 McAllester, D. A. (1999a). PAC-Bayesian model averaging. In *Proceedings of the twelfth annual*
1053 *conference on Computational learning theory*, pages 164–170. ACM.
- 1054 McAllester, D. A. (1999b). Some PAC-Bayesian theorems. *Machine Learning*, 37(3):355–363.
- 1055 Miller, J. W. and Dunson, D. B. (2019). Robust Bayesian inference via coarsening. *Journal of the*
1056 *American Statistical Association*, 114(527):1113–1125.
- 1057 Muandet, K., Fukumizu, K., Sriperumbudur, B., Schölkopf, B., et al. (2017). Kernel mean embed-
1058 ding of distributions: A review and beyond. *Foundations and Trends® in Machine Learning*,
1059 10(1-2):1–141.
- 1060 Neal, R. M. (2012). *Bayesian learning for neural networks*, volume 118. Springer Science &
1061 Business Media.

- 1062 Ovadia, Y., Fertig, E., Ren, J., Nado, Z., Sculley, D., Nowozin, S., Dillon, J., Lakshminarayanan, B.,
1063 and Snoek, J. (2019). Can you trust your model’s uncertainty? evaluating predictive uncertainty
1064 under dataset shift. *Advances in neural information processing systems*, 32.
- 1065 Polyanskiy, Y. and Wu, Y. (2014). Lecture notes on information theory. *Lecture Notes for ECE563*
1066 *(UIUC) and*, 6(2012-2016):7.
- 1067 Roberts, G. O. and Tweedie, R. L. (1996). Exponential convergence of Langevin distributions and
1068 their discrete approximations. *Bernoulli*, pages 341–363.
- 1069 Santambrogio, F. (2015). Optimal transport for applied mathematicians. *Birkäuser, NY*, 55.
- 1070 Santambrogio, F. (2017). {Euclidean, metric, and Wasserstein} gradient flows: an overview. *Bulletin*
1071 *of Mathematical Sciences*, 7:87–154.
- 1072 Shawe-Taylor, J. and Williamson, R. C. (1997). A PAC analysis of a Bayesian estimator. In *Annual*
1073 *Workshop on Computational Learning Theory: Proceedings of the tenth annual conference on*
1074 *Computational learning theory*, volume 6, pages 2–9.
- 1075 Veretennikov, A. Y. (2006). On ergodic measures for McKean-Vlasov stochastic equations. In *Monte*
1076 *Carlo and Quasi-Monte Carlo Methods 2004*, pages 471–486. Springer Berlin Heidelberg.
- 1077 Villani, C. (2003). *Topics in optimal transportation*, volume 58. American Mathematical Soc.
- 1078 Villani, C. et al. (2009). *Optimal transport: old and new*, volume 338. Springer.
- 1079 Welling, M. and Teh, Y. W. (2011). Bayesian learning via stochastic gradient Langevin dynamics.
1080 In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages 681–
1081 688.
- 1082 Wild, V. D., Hu, R., and Sejdinovic, D. (2022). Generalized variational inference in function spaces:
1083 Gaussian measures meet Bayesian deep learning. *Advances in Neural Information Processing*
1084 *Systems*, 35:3716–3730.
- 1085 Wilson, A. G. (2020). The case for Bayesian deep learning. *arXiv preprint arXiv:2001.10995*.
- 1086 Wilson, A. G. and Izmailov, P. (2020). Bayesian deep learning and a probabilistic perspective of
1087 generalization. *Advances in neural information processing systems*, 33:4697–4708.
- 1088 Wu, P.-S. and Martin, R. (2023). A comparison of learning rate selection methods in generalized
1089 Bayesian inference. *Bayesian Analysis*, 18(1):105–132.
- 1090 Zou, H. and Hastie, T. (2005). Regularization and variable selection via the elastic net. *Journal of*
1091 *the royal statistical society: series B (statistical methodology)*, 67(2):301–320.