FINITE-SAMPLE BOUNDS FOR ADAPTIVE INVERSE REINFORCEMENT LEARNING USING PASSIVE LANGEVIN DYNAMICS*

LUKE SNOW AND VIKRAM KRISHNAMURTHY [†]

Abstract. This paper provides a finite-sample analysis of a passive stochastic gradient Langevin dynamics algorithm (PSGLD) designed to achieve adaptive inverse reinforcement learning (IRL). By passive, we mean that the noisy gradients available to the PSGLD algorithm (inverse learning process) are evaluated at randomly chosen points by an external stochastic gradient algorithm (forward learner) that aims to optimize a cost function. The PSGLD algorithm acts as a randomized sampler to achieve adaptive IRL by reconstructing this cost function nonparametrically from the stationary measure of a Langevin diffusion. Previous work has analyzed the asymptotic performance of this passive algorithm using weak convergence techniques. This paper analyzes the non-asymptotic (finite-sample) performance using a logarithmic-Sobolev inequality and the Otto-Villani Theorem. We obtain finite-sample bounds on the 2-Wasserstein distance between the estimates generated by the PSGLD algorithm and the cost function. Apart from achieving finite-sample guarantees for adaptive ive IRL, this work extends a line of research in analysis of passive stochastic gradient algorithms to the finite-sample regime for Langevin dynamics.

Key words. stochastic gradient Langevin dynamics, passive learning, inverse reinforcement learning, finite-sample analysis, logarithmic-Sobolev inequality, Wasserstein distance

1. Introduction. We derive non-asymptotic bounds for a Langevin dynamics algorithm performing real-time inverse reinforcement learning (IRL). Traditional IRL [21], [12], [4] reconstructs the cost function of a Markov Decision Process by observing decisions taken from an optimal policy, i.e., *after* an observed agent has completed learning the optimal policy. Here, we consider *real-time* (adaptive) IRL. We observe an agent (forward learner) performing stochastic gradient descent (e.g., policy gradient reinforcement learning) on a cost function J, and attempt to reconstruct J in real-time. Thus, this technique can be regarded as an inverse stochastic gradient algorithm. It applies to IRL problems in several contexts such as adaptive Bayesian learning, constrained Markov Decision Processes, and logistic regression classification [16].

To accomplish real-time IRL, we employ a passive stochastic gradient Langevin dynamics (PSGLD) algorithm, initially proposed in [16]. Given observations of sequential stochastic gradient descent (SGD) evaluations on J, the PSGLD algorithm acts as a Markov chain Monte Carlo (MCMC) sampler designed to reconstruct J. Classical stochastic gradient Langevin dynamics [28], [10], evolves in the direction of a stochastic gradient evaluation (of a cost function) at the current iterate, plus an independent Gaussian perturbation. It asymptotically samples from the Gibbs measure encoding the cost function, and thus can be used to sample from probability distributions in the context of e.g., Bayesian learning [28] or empirical risk minimization [22]. Our PSGLD algorithm is considered *passive* because the stochastic gradient evaluations are not directly controlled, but are provided by the observed SGD process. Remarkably, [16] use stochastic approximation techniques to show that the PSGLD

^{*}The results in this paper have appeared in the Proceedings of the 2023 IEEE Conference on Decision and Control in reduced and abbreviated form. This manuscript substantially expands on the exposition of the Conference version, provides all proofs, and includes more detailed examples, discussion, and mathematical details.

Funding: This research was supported in part by the National Science Foundation grants CCF-2112457 and CCF-2312198, Army Research Office grant W911NF-19-1-0365, and Air Force Office of Scientific Research grant FA9550-22-1-0016

[†]Department of Electrical & Computer Engineering, Cornell University, Ithaca, NY

algorithm asymptotically samples from the Gibbs measure encoding the cost function being optimized by the observed SGD process. So the PSGLD algorithm asymptotically achieves inverse reinforcement learning. Similar passive schemes and stochastic approximation analyses have been investigated in [17], [30], [19].

In any practical IRL implementation it is necessary to understand how well the PSGLD algorithm recovers this cost function after a *finite run-time*. In this work we present a *non-asymptotic* analysis of this PSGLD algorithm; we provide finite-sample bounds on the 2-Wasserstein distance between the law of the algorithm and that of the Gibbs measure encoding the cost function. Non-asymptotic analysis of stochastic gradient Langevin dynamics has been investigated in [22], [8], [6]. In our case the algorithm is *passive*; so our analysis generalizes and extends previous works to handle this complexity.

Main Result: Recall J denotes the cost function being optimized by the forward learner, and which we aim to reconstruct. Denote π_k the sampling measure of our PSGLD algorithm at iterate $k \in \mathbb{N}$, π_{∞} the Gibbs measure w.r.t. $J(\pi_{\infty} \propto \exp(-\beta J))$, with β a controllable PSGLD parameter), and $\mathcal{W}_2(\pi_k, \pi_{\infty})$ the 2-Wasserstein distance between these. Observe that sampling from π_{∞} , achieved when $\mathcal{W}_2(\pi_k, \pi_{\infty}) = 0$, suffices to reconstruct J by taking the log-sample density. We obtain a bound on $\mathcal{W}_2(\pi_k, \pi_{\infty})$ that scales as $\mathcal{O}(k\epsilon\sqrt{\epsilon} + \exp(-k\epsilon))$, where ϵ is the algorithm step size. So we can choose $k\epsilon$ fixed and large enough to diminish $\exp(-k\epsilon)$, then ϵ small enough (with k simultaneously increasing to fix $k\epsilon$) to diminish $k\epsilon\sqrt{\epsilon}$ arbitrarily. Thus, for any y > 0 we can choose the step size ϵ small enough and algorithmic iterate k large enough such that $\mathcal{W}_2(\pi_k, \pi_{\infty}) \leq \mathcal{O}(y)$. The main result of this paper is a precise formulation of this statement.

Proof Technique: We bound $\mathcal{W}_2(\pi_k, \pi_\infty) \leq \mathcal{W}_2(\pi_k, \nu_{k\epsilon}) + \mathcal{W}_2(\nu_{k\epsilon}, \pi_\infty)$, where $\nu_{k\epsilon}$ is the measure, at time $k\epsilon$, of a particular continuous time diffusion with stationary measure π_{∞} . We obtain $\mathcal{W}_2(\pi_k,\nu_{k\epsilon}) \leq \mathcal{O}(k\epsilon\sqrt{\epsilon})$ through a Girsanov change of measure technique and a weighted transportation cost inequality. We then show that the diffusion satisfies a logarithmic-Sobolev inequality, allowing us to employ exponential decay of entropy and the Otto-Villani Theorem to show exponential decay of $\mathcal{W}_2(\nu_{k\epsilon}, \pi_{\infty})$. This proof structure is mirrored in [22]. However, our algorithm necessitates a non-trivial extension of the methods in [22]; we utilize both a generalized stochastic gradient Langevin dynamics form and a weighting kernel to control the external gradient evaluations. The generalized form disrupts absolute continuity of measure between the algorithm and continuous time diffusion, necessitating the introduction and control of an intermediate process in order to apply Girsanov's Theorem. It also necessitates control of the "sampling distribution" (from which initial SGD and PSGLD points are taken) to decrease discretization error. However, this control simultaneouly increases a relative entropy term appearing in the final 2-Wasserstein bound; this is handled by a careful specification of other algorithmic parameters. Finally, the continuous time diffusion is also distinct from that in [22], so we must prove logarithmic-Sobolev inequality satisfaction by a novel Lyapunov function. Furthermore, there are a host of supporting Lemmas, such as exponential integrability of the generalized diffusion, which have been necessarily derived for our analysis.

Another motivation for our work is the extension of a line of research in analysis of passive stochastic gradient algorithms. Such work [23] [13], [20], [31] has historically focused on passive stochastic approximation for e.g., sequential non-parametric estimation of regression functions. Recently, [16] extended this analysis to stochastic gradient Langevin dynamics in the passive regime. This paper serves as the first *finite-sample* result for stochastic gradient Langevin dynamics in the passive regime. **1.1. Organization.** Section 2 provides background on passive stochastic gradient Langevin dynamics and the asymptotic result of [16]. Section 3 discusses our main result, a non-asymptotic 2-Wasserstein bound (Theorem 3.9), and provides an example of how this IRL result interfaces with a canonical reinforcement learning algorithm. In section 4 we provide additional mathematical background for the proof of our bound, and section 5 provides details on the structure of our proof of Theorem 3.9. The complete proof details appear in the supplementary material.

2. Background: Passive Langevin Dynamics. Motivating our finite-sample analysis, this section introduces Langevin dynamics, presents the adaptive IRL setting and discusses the weak convergence asymptotic analysis of the PSGLD algorithm in [16].

2.1. Langevin Dynamics. The classical Langevin dynamics algorithm is given, with step size ϵ_k , objective function J, noise parameter β , and i.i.d. standard N-variate Gaussian noise w_k , as

(2.1)
$$\theta_{k+1} = \theta_k - \epsilon_k \nabla J(\theta_k) + \sqrt{2\epsilon\beta^{-1}} w_k, \quad k \in \mathbb{N}$$

Here $\theta_k \in \mathbb{R}^N$ is initialized by $\theta_0 \sim \pi_0$ for some sampling distribution π_0 on \mathbb{R}^N . The algorithm (2.1) is used primarily for either non-convex optimization [29] or to sample from probability distributions via MCMC [28]. The former is accomplished by treating (2.1) as a simulated annealer, and letting the step size ϵ_k and 'temperature' β^{-1} decrease to zero as $k \to \infty$. To accomplish the latter, the step size ϵ_k and temperature β^{-1} are *fixed* for all k. This work considers the latter case of constant step-size Langevin dynamics, with $\epsilon_k = \epsilon \ \forall k \in \mathbb{N}$. It is well known that the Markov process (2.1), with constant step-size, asymptotically samples from the *Gibbs* measure

(2.2)
$$\pi_{\infty}(\theta) = \frac{1}{\Lambda} \exp(-\beta J(\theta))$$

Here $\Lambda = \int_{\mathbb{R}^N} \exp(-\beta J(\alpha)) d\alpha$ is a normalizing constant and J is called the potential function. Indeed, (2.1) corresponds to a discretization of the continuous-time Langevin diffusion given by, with $\theta(t) \in \mathbb{R}^N$ and W(t) standard Brownian motion in \mathbb{R}^N , the Itô stochastic differential equation (SDE)

(2.3)
$$d\theta(t) = -\nabla J(\theta(t))dt + \sqrt{2\beta^{-1}dW(t)}, \quad t \ge 0$$

Under suitable conditions on J and β , this SDE has the Gibbs measure (2.2) as its unique stationary measure [7]. Thus, the Langevin dynamics algorithm (2.1) can be used as a MCMC algorithm to asymptotically sample from any probability distribution which can be expressed as (2.2) with some potential function J.

In [25] more general reversible diffusions of the form, with $\sigma : \mathbb{R}^N \to \mathbb{R}$ differentiable,

(2.4)
$$d\theta(t) = \left[-\frac{\beta}{2}\nabla J(\theta)dt - \nabla\sigma(\theta)dt + dW(t)\right]\sigma(\theta)$$

are studied, and it is shown that (2.4) has the same stationary measure (2.2) as the classical Langevin diffusion (2.3). The corresponding Euler-Maruyama time discretization results in the following discrete-time Markov process

(2.5)
$$\theta_{k+1} = \theta_k - \epsilon \left[\frac{\beta}{2}\nabla J(\theta_k) + \nabla \sigma(\theta_k)\right] \sigma(\theta_k) + \sqrt{\epsilon}\sigma(\theta_k)w_k$$

which can thus equivalently be used as a MCMC sampler from (2.2). We will utilize the generalized process (2.5), as opposed to the classical Langevin dynamics (2.1), for our PSGLD algorithm¹.

Active vs. Passive Gradient Evaluation: Notice that the above SGLD algorithms utilize, at each time step, the gradient $\nabla J(\theta_k)$ evaluated at the current iterate θ_k . We term this active gradient evaluation, and distinguish this from passive gradient evaluation, where the gradient is evaluated at a different (uncontrolled) point. The following section introduces the adaptive IRL setting and motivates the need for passive gradient evaluation in our PSGLD algorithm.

2.2. Adaptive IRL: Forward and Inverse Learning. This section motivates the setting of adaptive IRL. An agent (forward learner) is in the process of optimizing a cost function. An inverse learner passively observes sequential algorithmic iterates of the forward learner, and attempts to reconstruct (learn) the cost function being optimized. This paper provides non-asymptotic guarantees on the inverse learner's cost function reconstruction. Here we first introduce the dynamics of the forward learner, then the inverse learning setting and PSGLD algorithm.

2.2.1. Forward Learner. The forward learner runs a stochastic gradient descent (SGD) on a non-negative cost function

where the initial point evaluation is sampled randomly from sampling distribution $\pi_{0,\gamma}$ on \mathbb{R}^N . Examples include policy gradient reinforcement learning [26], neural network optimization [24] and federated learning [14]. We define the sampling distribution $\pi_{0,\gamma}$ as follows

(2.7)
$$\pi_{0,\gamma}(x) := \frac{\pi_0(\frac{x}{\gamma})}{\int_{\mathbb{R}^N} \pi_0(\frac{x}{\gamma}) dx}$$

where π_0 is an arbitrary density function and γ is a scale parameter ².

We assume the stochastic gradient algorithm (forward learner) resets after some finite time, so that the SGD process repeats indefinitely. This can be motivated by optimization of a non-convex function, in which re-initialization allows sufficient exploration, or multiple agents each learning to optimize the same cost function. Thus we have, for $n \in \mathbb{N}$ representing each "run" of the SGD, and τ_n stopping times:

(2.8)
$$\theta_{k+1} = \theta_k - \eta \nabla J(\theta_k), \quad k \in \{\tau_n, \dots, \tau_{n+1} - 1\}$$

where each $\theta_{\tau_n} \stackrel{i.i.d.}{\sim} \pi_{0,\gamma}$ and $\eta > 0$ is a fixed step-size. Here $\hat{\nabla}J(\theta_k)$ is an unbiased estimate of the true gradient $\nabla J(\theta_k)$, with bounded variance, see 3.2. Algorithm 2.1 displays this randomly re-initializing stochastic gradient descent.

Remark: The forward learner is not *necessarily* restricted to implementing an SGD algorithm; any process which provides sequential stochastic gradients of the cost function J, including i.i.d. samples from measure $\pi_{0,\gamma}$ or sequential algorithmic iterates with more sophisticated dependencies than (2.8) will suffice, see [16].

¹see [16] for motivation

²This construction is purely for notational convenience.

Algorithm 2.1 Randomly Re-Initializing SGD Process

 $\begin{array}{l} \text{initialize } \tau_0 = n = 0 \\ \textbf{while } n \geq 0 \ \textbf{do} \\ & \text{generate } \tau_{n+1} > \tau_n, \quad \theta_{\tau_n} \overset{iid}{\sim} \pi_{0,\gamma} \\ \textbf{for } k = \tau_n : \tau_{n+1} - 1 \ \textbf{do} \\ & \theta_{k+1} \leftarrow \theta_k - \eta \hat{\nabla} J(\theta_k) \\ & \textbf{end for} \\ \textbf{end while} \end{array}$

2.2.2. Inverse Learning: PSGLD. In this paper we take the perspective of an inverse learner who observes the SGD process (2.8), and attempts to reconstruct the cost function J being optimized. We assume this observer knows the sampling distribution $\pi_{0,\gamma}$ and can observe evaluations $\theta_k, k \in \mathbb{N}$. The agent recovers noisy gradient evaluations $\hat{\nabla} J(\theta_k) = \frac{\theta_{k+1} - \theta_k}{\eta}$.

Using only these sequential noisy gradient evaluations, how can the agent learn J? This is accomplished via MCMC sampling, using the following *passive stochastic gradient Langevin dynamics (PSGLD)* updates:

(2.9)
$$\alpha_{k+1} = \alpha_k - \epsilon \left[K_\Delta(\theta_k - \alpha_k) \frac{\beta}{2} \hat{\nabla} J(\theta_k) + \nabla \pi_{0,\gamma}(\alpha_k) \right] \pi_{0,\gamma}(\alpha_k) + \sqrt{\epsilon} \pi_{0,\gamma}(\alpha_k) w_k$$
$$\alpha_0 \sim \pi_{0,\gamma}$$

Note that α_0 is sampled randomly from the sampling distribution $\pi_{0,\gamma}$ of the SGD process (2.8). Here $\{w_k, k \ge 0\}$ is an i.i.d. sequence of standard *N*-variate Gaussian random variables,

(2.10)
$$K_{\Delta}(\theta_k - \alpha_k) := \frac{1}{\Delta^N} K\left(\frac{\theta_k - \alpha_k}{\Delta}\right)$$

is the Δ -parametrized kernel function, and β is the inverse temperature parameter. The algorithm is *passive* since the stochastic gradients $\hat{\nabla}J(\theta_k)$ and evaluation points θ_k are passively observed from SGD process (2.8). The kernel³ function $K(\cdot)$ controls for bias in these passive gradient evaluations, and can be chosen by the observer as any function $K : \mathbb{R}^N \to \mathbb{R}$ satisfying:

(2.11)
$$K(u) \ge 0, \quad K(u) = K(-u), \quad \sup_{u} K(u) < \infty,$$
$$\int_{\mathbb{R}^{N}} K(u) du = 1, \quad \int_{\mathbb{R}^{N}} |u|^{2} K(u) < \infty$$

 K_{Δ} weights the relevance of stochastic gradient $\hat{\nabla}J(\theta_k)$ to the current iterate $\alpha(t)$. We obtain K_{Δ} by modulating K by the domain scaling parameter Δ as (2.10). So Δ modulates the degree to which samples θ_k at a fixed distance from current iterate $\alpha(t)$ impact the algorithm's evolution.

Algorithm 2.2 displays this passive stochastic gradient Langevin dynamics algorithm, which takes as input the sequential evaluations θ_k made in Algorithm 2.1.

We claim that Algorithm 2.2 achieves adaptive inverse reinforcement learning, reconstructing J by taking only $\hat{\nabla}J(\theta_k)$ as input. Next a background result is presented which establishes this claim rigorously.

³An example kernel function is the multivariate normal $\mathcal{N}(0, \sigma^2 I_N)$ density with $\sigma = \Delta$, i.e., $\frac{1}{\Delta^N} K(\frac{\theta - \alpha}{\Delta}) = (2\pi)^{-N/2} \Delta^{-N} \exp(-\frac{\|\theta - \alpha\|^2}{2\Delta^2})$

Algorithm 2.2 PSGLD

parameters: step size ϵ , inverse temperature β , kernel scale Δ , re-sampling distribution scale γ initialize $\alpha_0 \sim \pi_{0,\gamma}$ while $k \ge 0$ do obtain θ_k from a Algorithm 2.1 if $k \ge 1$ then $\hat{\nabla}J(\theta_k) = \frac{1}{\eta}(\theta_k - \theta_{k-1}), \quad K_{k-1} = \frac{1}{\Delta^N}K(\frac{\theta_{k-1}-\alpha_{k-1}}{\Delta})$ sample $w_k \sim \mathcal{N}(0, I_N)$ $\alpha_k \leftarrow \alpha_{k-1} - \epsilon \left[K_{k-1}\frac{\beta}{2}\hat{\nabla}J(\theta_k) + \nabla\pi_{0,\gamma}(\alpha_{k-1})\right]\pi_{0,\gamma}(\alpha_{k-1}) + \sqrt{\epsilon}\pi_{0,\gamma}(\alpha_{k-1})w_k$ end if end while

2.3. Passive SGLD: Asymptotic Convergence. [16] provides the following result:

PROPOSITION 2.1 (Weak Convergence [16]). Let $\alpha^{\epsilon}(t) = \alpha_k$ for $t \in [\epsilon k, \epsilon(k+1)]$ be the continuous-time interpolation of PSGLD (2.9). Under assumptions (A1)-(A4) of [16], the process $\alpha^{\epsilon}(t)$ converges weakly to the solution of the stochastic differential equation

(2.12)

$$d\alpha(t) = -\left[\frac{\beta}{2}\pi_{0,\gamma}^2(\alpha(t))\nabla J(\alpha(t)) + \nabla\pi_{0,\gamma}(\alpha(t))\pi_{0,\gamma}(\alpha(t))\right]dt + \pi_{0,\gamma}(\alpha(t))dW(t)$$

$$\alpha(0) = \alpha_0 \sim \pi_{0,\gamma}$$

where W(t) is standard N-dimensional Brownian motion. Furthermore, the stochastic differential equation (2.12) has π_{∞} (2.2) as its stationary distribution.

Thus, we can use the algorithm (2.9) to generate asymptotic samples

$$\alpha_k \sim \pi_\infty(\alpha) \propto \exp(-\beta J(\alpha))$$

and reconstruct J from the sample log-density.

Motivation: Proposition 2.1 shows that Algorithm 2.2 asymptotically produces samples $\alpha_k \sim \pi_{\infty}$, and so the cost function J can be reconstructed from the logarithm of the asymptotic sample density. However, in this paper we are interested in quantifying how well this sampling algorithm approximates the Gibbs measure after a *finite run-time*. Our main result gives non-asymptotic (finite-sample) bounds on the 2-Wasserstein distance between the distribution of the sampling density produced by Algorithm 2.2 and the Gibbs measure π_{∞} (2.2).

3. Main Result. Non-Asymptotic Analysis of Passive Stochastic Gradient Langevin Dynamics. Recall from the Introduction that π_k is the sampling measure of Algorithm 2.2 at iterate k, π_{∞} is the Gibbs measure proportional to $\exp(-\beta J)$, and $\mathcal{W}_2(\pi_k, \pi_{\infty})$ is the 2-Wasserstein distance between these. Our main result is as follows: for any y > 0 we can choose the step size ϵ small enough and iteration number k large enough such that $\mathcal{W}_2(\pi_k, \pi_{\infty}) \leq \mathcal{O}(y)$. In this section we formulate this result precisely. We provide a brief overview of the 2-Wasserstein metric, specify assumptions on the cost function J and sampling distribution $\pi_{0,\gamma}$, provide the main bound in the form of Theorem 3.9 , and discuss the application to adaptive inverse reinforcement learning in a Markov Decision Process framework. **3.1. 2-Wasserstein Distance.** We provide a non-asymptotic bound on the convergence of (2.9) to the Gibbs measure π_{∞} (2.2), in terms of the 2-Wasserstein distance:

(3.1)
$$\mathcal{W}_2(\mu,\nu) := \inf_{\gamma \in \Gamma(\mu,\nu)} \left(\mathbb{E}_{(x,y) \sim \gamma} \|x-y\|^2 \right)^{1/2}$$

Here $\Gamma(\mu, \nu)$ is the set of all couplings of measures μ and ν , where a coupling γ is a joint probability measure on $\mathbb{R}^N \times \mathbb{R}^N$ with marginals μ and ν , i.e., $\gamma(A, \mathbb{R}^N) = \mu(A), \gamma(\mathbb{R}^N, B) = \nu(B), \forall A, B \in \mathcal{B}(\mathbb{R}^N)$, where $\mathcal{B}(\mathbb{R}^N)$ is the Borel σ -algebra of \mathbb{R}^N . Notice that the Wasserstein distance (3.1) indeed satisfies all axioms of a metric on the space of measures. The 2-Wasserstein distance is a more suitable metric for assessing the quality of approximate sampling schemes [8], [22], than others such as total-variation norm, since it gives direct guarantees on the accuracy of approximating higher order moments [8]. However, it also precludes us from using SDE discretization analysis presented in the seminal book [?], which utilizes total-variational norm.

3.2. Assumptions. Here we list several assumptions on the cost function J (2.6) of the forward learner and the base sampling distribution π_0 , required for the finite-sample analysis. Assumptions on J are standard and equivalent to those taken in [22]. Assumptions on the base sampling distribution π_0 hold for a wide class of probability density functions, including Gaussian densities.

A 3.1 (J regularity). J is L_J -Lipschitz continuous and $L_{\nabla J}$ -smooth: $\exists L_J$, $L_{\nabla J} > 0$ such that for all $x, y \in \mathbb{R}^N$,

$$||J(x) - J(y)|| \le L_J ||x - y||, \qquad ||\nabla J(x) - \nabla J(y)|| \le L_{\nabla J} ||x - y||$$

A 3.2 (Dissipativity). J is (m, b)-dissipative:

$$\exists m > 0, b \ge 0 : \langle x, \nabla J(x) \rangle \ge m \|x\|^2 - b, \ \forall x \in \mathbb{R}^N$$

A 3.3 (Gradient Noise Variance). The noisy SGD gradient evaluation is unbiased, i.e. $\mathbb{E}[\hat{\nabla}J(x)] = \nabla J(x) \ \forall x \in \mathbb{R}^N$. Furthermore, the noise is additive such that $\hat{\nabla}J(x) - \nabla J(x)$ is i.i.d. with variance bounded uniformly in x, i.e., there exists a constant $\zeta \geq 0$ such that

$$\mathbb{E}[\|\hat{\nabla}J(x) - \nabla J(x)\|^2] \le \zeta, \ \forall x \in \mathbb{R}^N$$

A 3.4 (π_0 Exponential Decay). The base sampling distribution π_0 has an exponential tail decay and differential decay $\mathcal{O}(||x||^{-1})$, i.e.,

$$\exists M \in \mathbb{N}, \tilde{C} > 0 : \pi_0(x) \le \exp(-\|x\|^2), \quad \|\nabla \pi_0(x)\| \le \frac{C}{\|x\|} \quad \forall \|x\| > M$$

A 3.5 (π_0 Lipschitz-continuity).

$$\exists L_{\pi_0} > 0 : \|\pi_0(x) - \pi_0(y)\| \le L_{\pi_0} \|x - y\| \ \forall x, y \in \mathbb{R}^N$$

A 3.6 (π_0 Structure). π_0 is unimodal and has support on \mathbb{R}^N .

A 3.7 (Kernel Structure). The kernel function $K(\cdot)$ satisfies (2.11).

A 3.8 (Feasible Parameter Ranges). Here \wedge denotes the min operator and \vee the max operator. Assume

$$i) \ \eta \in (0, 1 \land \frac{m}{4L_{\nabla J}^2})$$

$$ii) \ \epsilon \in \left(0, 1 \land \sqrt{\frac{1}{249}L_{\nabla J}^{-1}}\right)$$

$$iii) \ \beta \ge \frac{1}{4L_{\nabla J}^2} \lor \frac{\sqrt{2\pi+4}}{m\sqrt{L_{\nabla J}}}$$

7



Figure 1: High level procedure for achieving inverse reinforcement learning. The forward learning process is represented by a stochastic gradient descent (SGD), and the inverse learner incorporates sequential SGD evaluations θ_k into its PSGLD algorithm to reconstruct J. The PSGLD algorithm reconstructs J by approximately sampling from the Gibbs measure π_{∞} (then taking the log-sample density). We measure the proximity of the PSGLD algorithm to π_{∞} by $W_2(\pi_k, \pi_{\infty})$, the 2-Wasserstein distance between the sample law of α_k and the measure π_{∞} . We control this distance by bounding it by $W_2(\pi_k, \nu_{k\epsilon}) + W_2(\nu_{k\epsilon}, \pi_{\infty})$, where $\nu_{k\epsilon}$ is the law of $\alpha(t)$ at time $t = k\epsilon$.

3.2.1. Discussion of assumptions. A3.1 - A3.4 are equivalent to those taken for the objective function in [22]. A3.1 is widely used in the literature on non-convex optimization and sampling. A3.2 can be enforced through weight decay regularization [18], see section 4 of [22] for more details. A3.3 is a standard assumption for stochastic gradient evaluations. A3.4 - A3.6 admit a wide range of probability density functions, including Gaussians. We note that for A3.8 to be satisfied in practice, the inverse learner must have some knowledge of feasible ranges for Hessian bound $L_{\nabla J}$ and dissipativity constant m; once these ranges are known then ϵ can be taken small enough and β large enough so that (*ii*) and (*iii*) are satisfied. Notice that the feasible range for η can always be satisfied; the SGD process (2.8) optimizing cost function Jwith step $\hat{\eta} \geq (1 \wedge \frac{m}{4L_{\nabla J}^2})$ is equivalent to another SGD with step $\eta < \frac{m}{4L_{\nabla J}^2}$ which optimizes $\frac{\eta}{\hat{\eta}}J$. So assuming η which satisfies A3.8 we can sample from $\pi_{\infty} \propto \exp(-\frac{\eta}{\hat{\eta}}\beta J)$, from which J can be recovered since the scale $\frac{\eta}{\hat{\eta}}\beta$ disappears upon MCMC sample measure normalization.

3.3. Main Result. Finite-Sample Bound. Letting

$$\pi_k := \operatorname{Law}(\alpha_k), \quad \nu_{k\epsilon} := \operatorname{Law}(\alpha(k\epsilon))$$

be the respective measures of the sampling density produced by iterates α_k (2.9) and the continuous time diffusion $\alpha(t)$ (2.12) at time $t = k\epsilon$, we may bound

$$\mathcal{W}_2(\pi_k, \pi_\infty) \le \mathcal{W}_2(\pi_k, \nu_{k\epsilon}) + \mathcal{W}_2(\nu_{k\epsilon}, \pi_\infty)$$

Figure 1 shows the high level procedure for achieving inverse reinforcement learning. The forward learner is represented by a stochastic gradient descent (SGD) process which optimizes J. The PSGLD algorithm takes in sequential SGD evaluations θ_k and produces samples α_k which approximately sample from the Gibbs measure π_{∞} , allowing for reconstruction of J by taking the log-sample density. We measure this approximation by the distance $\mathcal{W}_2(\pi_k, \pi_{\infty})$, which can be bounded by introducing the intermediate continuous-time diffusion (2.12), since PSGLD (2.9) is an approximate discretization of (2.12) and (2.12) has Gibbs measure (2.2) as its stationary measure.

We present our Wasserstein bound in a way that explicitly depends on a hyperparameter δ , e.g., $W_2(\pi_k, \pi_\infty) \leq f(\delta)$ for some function f which is monotonically increasing and has $\lim_{\delta \to 0} f(\delta) = 0$. Both the Wasserstein bound and certain algorithmic parameters have a functional dependence on δ : for decreasing δ (decreasing $W_2(\pi_k, \pi_\infty)$), we require e.g., increasing the algorithmic iterations and decreasing the step size. Specifically, our main result states that for any arbitrarily small $f(\delta)$, we can choose the step size ϵ small enough, algorithmic iterations k large enough, kernel scale parameter Δ small enough, and sampling distribution scale parameter γ small enough, such that $W_2(\pi_k, \pi_\infty) \leq f(\delta)$. Next these qualitative parameter specifications are shown explicitly, as functions of control hyperparameter δ . Then our main bound on $W_2(\pi_k, \pi_\infty)$ is presented in Theorem 3.9.

3.3.1. Algorithmic Parameter Specifications. Here we show the dependence of algorithmic parameters on the hyperparameter δ , which controls the main Wasserstein bound presented in Theorem 3.9. δ acts as a one-dimensional "knob" that can be turned, which reveals the step size ϵ , iteration number k, etc., required to achieve a 2-Wasserstein bound proportional to δ . The main idea is that Theorem 3.9 presents a (monotonically increasing) function $f(\delta)$, with $\lim_{\delta \to 0} f(\delta) = 0$, such that for any $\delta > 0$ we can take algorithmic parameters as follows to obtain $W_2(\pi_k, \pi_\infty) \leq f(\delta)$.

Step Size:

(3.2)
$$\epsilon \le \left(\frac{\delta}{\log\left(\frac{1}{\delta}\right)}\right)^2 \wedge 1$$

Algorithmic Iterations:

(3.3)
$$k\epsilon \ge \beta c_{LS} \log\left(\frac{1}{\delta}\right)$$

where c_{LS} is the logarithmic-Sobolev constant of diffusion (2.12), explicitly bounded in (5.5).

Kernel Scale: Recalling, for general $\alpha \in \mathbb{R}_+$, $K_{\alpha}(\cdot) = \frac{1}{\alpha^N} K(\frac{\cdot}{\alpha})$, define $\hat{K}_{\alpha} := \sup_{x \in \mathbb{R}^N} K_{\alpha}(x)$. Also let K^{-1} denote the inverse of K and K^{-2} denote the inverse of K^2 , both mapping to the non-negative orthant, i.e., for $x \in \mathbb{R}$, $K^{-1}(x) := \{y \in \mathbb{R}^N_+ : K(y) = x\}$, $K^{-2}(x) := \{y \in \mathbb{R}^N_+ : K^2(y) = x\}$ where \mathbb{R}^N_+ is the set of N-dimensional vectors with all non-negative elements. This definition is without loss of generality, since K is chosen to be symmetric by (2.11). Then take

(3.4)
$$\Delta \leq \inf_{x \in [\epsilon, \hat{K}_{\epsilon}]} \frac{K^{-1}(\frac{\hat{K}_1 \sqrt{2\pi}}{2\epsilon} e^{x^2/2})}{K^{-2}(x \epsilon^{2N})}$$

Sampling Distribution Scale: Choose the base sampling distribution π_0 such that $\bar{\pi}_0 := \sup_x \pi_0(x) = 1$, and sampling distribution scale parameter γ as

(3.5)
$$\gamma \in [\epsilon^2, \epsilon^{3/2}]$$

3.3.2. Main 2-Wasserstein Bound.

THEOREM 3.9 (Finite-Sample 2-Wasserstein Bound). Consider the PSGLD Algorithm 2.2 with iterates $\alpha_k \in \mathbb{R}^N$. Recall c_{LS} is the logarithmic-Sobolev constant for Itô diffusion (2.12), bounded in (5.5). For any

(3.6)
$$\delta \in \left[0, \exp\left(-\frac{1}{\beta c_{LS}}\right)\right]$$

choose step size ϵ according to (3.2), number of iterations k according to (3.3), kernel scale Δ according to (3.4), and sampling distribution $\pi_{0,\gamma}$ with γ satisfying (3.5). Then, under assumptions (A3.1)-(A3.8), the 2-Wasserstein distance between the distribution π_k , generated by the PSGLD algorithm, and the Gibbs measure π_{∞} (2.2), satisfies:

(3.7)
$$\mathcal{W}_2(\pi_k, \pi_\infty) \le \delta \left[C_4 + \sqrt{2c_{LS}C_3} \right] + \delta \sqrt{10c_{LS}N \log(1/\delta)}$$

 C_3, C_4 are constants dependent on structural specifications of J and the process (2.9), provided explicitly in supplementary material 7.1. c_{LS} is the logarithmic-Sobolev constant bounded explicitly in Proposition 5.3.

Bound Discussion: For any $\alpha > 0$, $\delta \sqrt{\alpha \log(1/\delta)}$ is monotonically increasing in δ for $\delta \in (0, 0.607)$ and

$$\lim_{\delta \to 0} \delta \sqrt{\alpha \log \left(1/\delta \right)} = 0$$

So, $\mathcal{W}_2(\pi_k, \pi_\infty)$ is monotonically increasing in δ for $\delta \in (0, 0.607)$ and

$$\lim_{\delta \to 0} \mathcal{W}_2(\pi_k, \pi_\infty) = 0$$

Thus, Theorem 3.9 asserts that, through hyperparameter δ , we can control the number of iterations k as (3.3), step size ϵ as (3.2), kernel scale Δ as (3.4) and sampling distribution scale γ as (3.5), such that the PSGLD algorithm (2.9) is within any arbitrarily small desired 2-Wasserstein distance (3.7) to the Gibbs distribution (2.2). Here δ acts as a precision parameter; smaller δ yields a tighter approximation (3.7) at the expense of larger number of iterations k and smaller step size ϵ , kernel scale Δ and sampling distribution scale γ .

Recalling $\pi_{\infty}(\alpha) \propto \exp(-\beta J(\alpha))$, the cost function J can be approximately reconstructed as the logarithm of sample density produced by α_k . This reconstruction approaches the true cost function J as $\delta \to 0$. This result generalizes the nonasymptotic bound obtained in [22] (equation 3.3) to our *passive* stochastic gradient Langevin dynamics algorithm.

Parameter Specifications Discussion: Observe the parameter specifications (3.2) - (3.5) necessary for acheiving a given Wasserstein bound (3.7). Specifications (3.2) and (3.3) are intuitive; as we decrease the step size ϵ we should decrease the discretization error between algorithm (2.9) and continuous diffusion (2.12), and as we increase the iterations k we will decrease the distance from the diffusion (2.12) to its stationary measure π_{∞} .

However, the algorithm (2.9) is not an exact discretization of diffusion (2.12), as it has a gradient term governed by the external SGD process (2.8). The weighting kernel K is introduced to control for biases in this SGD-evaluated gradient, but for any non-zero variance of K there will still be biased gradient evaluations entering the algorithm (2.9) which prevent it from converging to π_{∞} . To minimize these, the kernel scale Δ can be reduced; however note that this should come at the cost of increasing the time needed to reach a specified Wasserstein bound, since "useful" gradient information will be integrated into the passive algorithm less often. We see this as an unavoidable tradeoff, one which necessitates taking Δ as (3.4), but which has not been fully quantified in this work.

Note that we also require decreasing the sampling distribution scale γ to obtain a tighter Wasserstein bound. This has arisen as a quantitative necessity in obtaining bounds in Lemma 5.1 and Lemma 5.2 (which are key developments in the proof of Theorem 3.9). The intuition is as follows: this specification allows us to control $\mathbb{E}\|\pi_{0,\gamma}(\alpha_k)\|^2$ and $\mathbb{E}\|\nabla\pi_{0,\gamma}(\alpha_k)\|^2$, such that the influence of $\pi_{0,\gamma}(\alpha_k)$ in the algorithm (2.9) does not outweigh that of $K_{\Delta}(\theta_k, \alpha_k)\hat{\nabla}J(\theta_k)$; (as in the previous paragraph explanation) as Δ gets smaller, $K_{\Delta}(\theta_k, \alpha_k)\hat{\nabla}J(\theta_k)$ contributes less often (but is more accurate when it does), and the contribution of $\pi_{0,\gamma}(\alpha_k)$ should balance this. Notice that we also must choose the base distribution π_0 wide enough (such that $\bar{\pi}_0 = 1$), and scale parameter not *too* small (lower bounded by ϵ^2), so that there is always some non-zero probability of sampling from any point in the domain, allowing for sufficient exploration.

3.4. Example. Adaptive MDP Inverse Reinforcement Learning. As discussed in [16], the PSGLD algorithm can be applied to a variety of settings including constrained MDPS, adaptive Bayesian inference, and logistic regression classification. Here we illustrate how the PSGLD algorithm (Alg 2.2), and finite-sample guarantees of Theorem 3.9, interface with a well-known practical reinforcement learning (RL) algorithm. We first outline the details of the forward RL procedure, then show how our PSGLD algorithm can be used to achieve inverse RL with nonasymptotic accuracy guarantees.

3.4.1. Forward Reinforcement Learning: Policy Gradient. Here we present one canonical and well-established algorithm for reinforcement learning (RL), the Reinforce algorithm. Consider the following notation

- discount factor $\gamma \in (0, 1)$, discrete-time index $t \in \mathbb{N}$
- states $s_t \in \mathcal{S}$, actions $a_t \in \mathcal{A}$, cost $c_t = C(s_t, a_t), C : \mathcal{S} \times \mathcal{A} \to \mathbb{R}_+$
- state transition probabilities $P^a_{s,s^\prime} = P(s_{t+1}=s^\prime|s_t=s,a_t=a)$
- initial state probability distribution $\rho_0(s) = \mathbb{P}(s_0 = s)$

- policy function $\pi(s, a; \theta) = \mathbb{P}[a_t = a | s_t = s; \theta]$ under policy parameter $\theta \in \mathbb{R}^d$ Now let

$$J(\theta) := \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t c_t \right] : \mathbb{R}^d \to \mathbb{R}_+$$

be the expected cost objective, where the expectation is taken with respect to the probability distribution of trajectories $(s_0, a_0, \ldots, s_t, a_t, \ldots)$ induced by ρ_0 , $P^a_{(s,s')}$, and $\pi(s, a; \theta)$. The reinforcement learning goal is to find a policy parameter θ such that $J(\theta)$ is minimized. This minima can be achieved by sequentially updating θ via stochastic gradient descent on $J(\theta)$. This is the methodology of the canonical Reinforce algorithm, which has sequential policy parameter updates given by

(3.8)
$$\theta_{k+1} = \theta_k - \eta \sum_{t=0}^T \left[\gamma^t \nabla_\theta \log \pi(s_t, a_t; \theta) \sum_{k=t}^T \gamma^{k-t} r_k \right] = \theta_k - \eta \hat{\nabla} J(\theta)$$

where, crucially, $\hat{\nabla} J(\theta)$ is an unbiased estimate of $\nabla J(\theta)$ by the Policy Gradient Theorem [26].

3.4.2. Adaptive Inverse Reinforcement Learning: PSGLD. Observe that the RL algorithm 3.8 works by employing a stochastic gradient descent in the space of policies. This suggests that we can employ our PSGLD algorithm, taking as input the sequential policy evaluations and outputting the expected cost function.

Suppose we observe an agent performing policy gradient RL by enacting sequential sample paths $\{(s_0, a_0, \ldots, s_T, a_T)_k \sim \pi(\cdot, \cdot; \theta_k)\}_{k=1}^K$, and updating θ_k according to Algorithm 3.8. Notice that by observing a sample path $(s_0, a_0, \ldots, s_T, a_T) \sim \pi(\cdot, \cdot; \theta)$ we can obtain an unbiased estimate $\hat{\theta}$ of θ by simply taking the empirical distribution of observed state-action pairs. Then form $\tilde{\nabla}J(\theta_k) = \frac{\hat{\theta}_k - \hat{\theta}_{k-1}}{\eta}$, and notice that

$$\mathbb{E}\left[\tilde{\nabla}J(\theta_k)|\,\theta_k,\,\theta_{k-1}\right] = \hat{\nabla}J(\theta_k)$$

i.e., $\tilde{\nabla}J(\theta_k)$ is a second order stochastic gradient estimate; it is an unbiased estimate of the stochastic gradient $\hat{\nabla}J(\theta_k)$, and $\hat{\nabla}J(\theta_k)$ is fully determined by the (random) evaluations θ_k , θ_{k+1} . So we have

$$\mathbb{E}\left[\tilde{\nabla}J(\theta_k)\right] = \mathbb{E}\left[\mathbb{E}\left[\tilde{\nabla}J(\theta_k)|\,\theta_k,\,\theta_{k+1}\right]\right] = \nabla J(\theta_k)$$

and thus $\nabla J(\theta_k)$ can be utilized in our passive stochastic gradient Langevin dynamics algorithm (Alg 2.2) since it is an unbiased estimate of $\nabla J(\theta_k)$. Thus, Theorem 3.9 applies and we can control the non-asymptotic Wasserstein distance arbitrarily. Fixing any $\delta > 0$, Algorithm 2.2, with parameters specified by (3.2) and with θ_k replaced by $\hat{\theta}_k$ above, produces iterates α_l , $l \geq k$, that sample from a measure π with

$$\mathcal{W}_2(\pi.\pi_\infty) \le \mathcal{O}(\delta + \delta\sqrt{\log(1/\delta)}), \quad \pi_\infty(\theta) \propto \exp(-\beta J(\theta))$$

Then J can be approximately recovered by taking the logarithm of MCMC sample density.

Note that traditional IRL methods aim to reconstruct C(s, a), rather than $J(\theta)$, given *optimal* policy demonstrations. In our case C(s, a) can be recovered up to a constant multiplicative factor once $J(\theta)$ and the MDP transition dynamics are known, since $J(\theta)$ is the expectation of C(s, a) with respect to the stationary measure induced by the policy $\pi(\cdot, \cdot; \theta)$ and the dynamics $P_{s,s'}^a$. Furthermore, in contrast to traditional methods [32], [21], we operate in the transient regime where the observed agent is *in* the process of learning an optimal policy.

4. Preliminaries for Bound Proof. Since finite-sample bounds for stochastic diffusions may not be widely known in the stochastic control community, this short section briefly summarizes the main tools required for the proof of Theorem 3.9.

4.1. Infinitesimal Generator. Let X_t be an \mathbb{R}^N -valued diffusion defined by the stochastic differential equation

(4.1)
$$dX_t = b(X_t) dt + \sigma(X_t) dW(t), \quad X_0 = x \in \mathbb{R}^N$$

where $b : \mathbb{R}^N \to \mathbb{R}^N$ is the drift function, $\sigma : \mathbb{R}^N \to \mathbb{R}$ is the volatility function, and W(t) is standard N-dimensional Brownian motion. Fixing a point $x \in \mathbb{R}^N$, let P^x denote the law of X_t given $X_0 = x$, and \mathbb{R}^x denote expectation with respect to P^x . Let \mathcal{L} be the *infinitesimal generator* of X_t , defined by its action on compactly-supported C^2 functions $f : \mathbb{R}^N \to \mathbb{R}$ in domain $\mathcal{D}(\mathcal{L}) \subseteq C^2(\mathbb{R}^N)$, as

(4.2)
$$\mathcal{L}f(x) = \lim_{t \downarrow 0} \frac{\mathbb{E}^x[f(X_t) - f(x)]}{t} = \sum_{i=1}^n b_i(x) \frac{\partial f}{\partial x_i}(x) + \frac{1}{2} \sum_{i,j} \sigma^2(x) \frac{\partial^2 f}{\partial x_i \partial x_j}(x)$$

where $b_i(x)$ is the *i*'th element of $b(x) \in \mathbb{R}^N$. Thus \mathcal{L} is an operator acting on $f \in C^2(\mathbb{R}^N)$ as

$$\mathcal{L}f = \frac{1}{2}\sigma^2 \Delta f + \langle b, \nabla f \rangle$$

where $\Delta := \nabla \cdot \nabla$ denotes the standard Laplacian operator. We say π is an invariant probability measure w.r.t \mathcal{L} if and only if $\int_{\mathbb{R}^N} \mathcal{L}gd\pi = 0$ for all $g \in \mathcal{D}(\mathcal{L})$.

In this work we consider the diffusion which solves the stochastic differential equation (2.12), which has:

$$b(x) = -\frac{\beta}{2}\pi_{0,\gamma}^{2}(x)\nabla J(x) - \pi_{0,\gamma}(x)\nabla\pi_{0,\gamma}(x), \quad \sigma(x) = \pi_{0,\gamma}(x)$$

Thus, the infinitesimal generator of our diffusion process (2.12) is given as

(4.3)
$$\mathcal{L}f = \frac{1}{2}\pi_{0,\gamma}^2 \Delta f - \frac{\beta}{2}\pi_{0,\gamma}^2 \langle \nabla J, \nabla f \rangle - \pi_{0,\gamma} \langle \nabla \pi_{0,\gamma}, \nabla f \rangle$$

and note that by assumptions (A2), (A6) and by Theorem 2.5 of [15], we have that (2.12) admits a unique strong solution.

4.2. Poincaré and logarithmic Sobolev inequalities. Considering a general infinitesimal generator \mathcal{L} , with stationary measure π , we can define the Dirichlet form

$$\mathcal{E}(g) := -\int_{\mathbb{R}^N} g\mathcal{L}g d\pi$$

and the spectral gap λ as

(4.4)
$$\lambda := \inf \left\{ \frac{\int_{\mathbb{R}^N} \mathcal{E}(g) d\pi}{\int_{\mathbb{R}^N} g^2 d\pi} : g \in C^1(\mathbb{R}^N) \cap L^2(\pi), g \neq 0, \int_{\mathbb{R}^N} g d\pi_\infty = 0 \right\}$$

Let us consider a Markov process X_t with unique invariant distribution π and infinitesimal generator \mathcal{L} . We say that π satisfies a *Poincaré (spectral gap) inequality* with constant c if

(4.5)
$$\chi^2(\mu||\pi) \le c \mathcal{E}\left(\sqrt{\frac{d\mu}{d\pi}}\right)$$

for all probability measures $\mu \ll \pi$ (μ absolutely continuous w.r.t π), where $\chi^2(\mu||\pi) := ||\frac{d\mu}{d\pi} - 1||^2_{L^2(\pi)}$ is the χ^2 divergence between μ and π . If (4.5) is satisfied for some c, then we have $\frac{1}{c} \leq \lambda$ where λ is the spectral gap given in (4.4). In particular, letting c_P denote the *Poincaré constant*, given as the smallest c such that (4.5) holds,

$$c_P = \inf\{c: \chi^2(\mu||\pi) \le c \mathcal{E}\left(\sqrt{\frac{d\mu}{d\pi}}\right) \ \forall \mu \ll \pi\}$$

where \ll denotes absolute continuity, then we have $\frac{1}{c_P} = \lambda$, and the eigenspectrum of $-\mathcal{L}$ is contained in $\{0\} \cup [\frac{1}{c_P}, \infty)$. We say that π satisfies a *logarithmic Sobolev inequality* with constant c if

$$D(\mu||\pi) \le 2 c \mathcal{E}\left(\sqrt{\frac{d\mu}{d\pi}}\right)$$
13

for all $\mu \ll \pi$, where $D(\mu || \pi) = \int d\mu \log \frac{d\mu}{d\pi}$ is the Kullback-Leibler divergence. One of the main efforts of this work will be to show that the diffusion (2.12)

One of the main efforts of this work will be to show that the diffusion (2.12) satisfies a log-Sobolev inequality; this then allows us to utilize several useful properties in the non-asymptotic analysis. Specifically, letting $\{X(t)\}_{t\geq 0}$ be a Markov process with stationary distribution π and Dirichlet form \mathcal{E} , then we have:

LEMMA 4.1 (Exponential decay of entropy [2], Th. 5.2.1). Let $\mu_t := Law(X(t))$. If π satisfies a logarithmic-Sobolev inequality with constant c, then

(4.6)
$$D(\mu_t || \pi) \le D(\mu_0 || \pi) e^{-2t/c}$$

LEMMA 4.2 (Otto-Villani theorem [2], Th. 9.6.1). If π satisfies a logarithmic-Sobolev inequality with constant c, then, for any $\mu \ll \pi$

(4.7)
$$\mathcal{W}_2(\mu, \pi) \le \sqrt{2cD(\mu||\pi)}$$

The following two results give sufficient conditions for a measure π to satisfy Poincare and logarithmic-Sobolev inequalities, using Lyapunov function criteria.

PROPOSITION 4.3 (Bakry 2008 [1]). Let $\pi(dx) = \exp(-H(x))dx$ be a probability measure on \mathbb{R}^N with $H \in C^2(\mathbb{R}^N)$ and lower bounded. Let \mathcal{L} be the infinitesimal generator of a Markov process with stationary measure π . Suppose there exist constants $\kappa_0, \zeta_0 > 0, r \ge 0$ and a C^2 function $V : \mathbb{R}^N \to [1, \infty)$ such that

(4.8)
$$\frac{\mathcal{L}V(w)}{V(w)} \le -\zeta_0 + \kappa_0 \mathbf{1}\{\|w\| \le r\}$$

Then π satisfies a Poincaré inequality with constant

(4.9)
$$c_P \le \frac{1}{\zeta_0} \left(1 + C\kappa_0 r^2 \exp(O_r(H)) \right)$$

where C > 0 is a universal constant and $O_r(H) := \max_{\|w\| \le r} H(w) - \min_{\|w\| \le r} H(w)$

PROPOSITION 4.4 (Cattiaux et. al. (2010) [5]). Let $\pi(dx) = \exp(-H(x))dx$ be a probability measure on \mathbb{R}^N with $H \in C^2(\mathbb{R}^N)$ and lower bounded. Let \mathcal{L} be the infinitesimal generator of a Markov process with stationary measure π . Suppose the following conditions hold:

1. There exist constants $\kappa, \gamma > 0$ and a C^2 function $V : \mathbb{R}^d \to [1, \infty)$ such that

(4.10)
$$\frac{\mathcal{L}V(w)}{V(w)} \le \kappa - \gamma \|w\|^2 \ \forall w \in \mathbb{R}^d$$

2. π_{∞} satisfies a Poincaré inequality with constant c_P .

3. There exists some constant $K \ge 0$, such that $\nabla^2 H \succcurlyeq -KI_d$

Let Z_1, Z_2 be defined, for some $\zeta > 0$ as

(4.11)
$$Z_1 = \frac{2}{\gamma} \left(\frac{1}{\zeta} + \frac{K}{2} \right) + \zeta, \quad Z_2 = \frac{2}{\gamma} \left(\frac{1}{\zeta} + \frac{K}{2} \right) \left(\kappa + \gamma \int_{\mathbb{R}^N} \|w\|^2 \pi(dw) \right)$$

Then π satisfies a logarithmic Sobolev inequality with constant $c_{LS} = Z_1 + (Z_2 + 2)c_P$.

We will be interested in showing that the invariant (Gibbs) measure π_{∞} (2.2) of our particular diffusion (2.12) satisfies a log-Sobolev inequality, so that we can apply Lemmas 4.1 and 4.2 to obtain exponential convergence of $W_2(\nu_{k\epsilon}, \pi_{\infty})$. To show the log-Sobolev inequality holds, we will show that the conditions of Proposition 4.4 hold, using Proposition 4.3 as an intermediate step. More details on this procedure will be outlined in Section 5.

The following result is unrelated to Poincaré and log-Sobolev inequalities, but gives a way to bound a general Wasserstein distance once a KL-divergence is known. We will utilize this in the bound on $W_2(\pi_k, \nu_{k\epsilon})$.

COROLLARY 4.5 (Bolley and Villani 2005 [3] Cor. 2.3). For any two Borel probability measures μ, ν on \mathbb{R}^N ,

$$\mathcal{W}_{2}(\mu,\nu) \leq 2 \inf_{\lambda>0} \left(\frac{1}{\lambda} \left(\frac{3}{2} + \log \int_{\mathbb{R}^{N}} e^{\lambda \|w\|^{2}} \nu(dw) \right) \right)^{1/2} \left[\sqrt{D(\mu||\nu)} + \left(\frac{D(\mu||\nu)}{2} \right)^{1/4} \right]$$

5. Proof of Main Result (Theorem 3.9). Outline. Here we provide the proof structure for our bound on $W_2(\pi_k, \pi_\infty)$, provided as (3.7) in Theorem 3.9 (complete proofs can be found in the supplementary material). The block diagram in Figure 2 displays the relations between our main supporting results in the proof of Theorem 3.9. The high level proof structure is as follows: We bound $W_2(\pi_k, \pi_\infty) \leq W_2(\pi_k, \nu_{k\epsilon}) + W_2(\nu_{k\epsilon}, \pi_\infty)$, i.e., we first control the discretization error between passive algorithm 2.2 and diffusion 2.12, then control the convergence rate of this diffusion to its stationary distribution π_∞ .

In order to achieve a useful bound on the former, scaling as $\mathcal{O}(k\epsilon\sqrt{\epsilon})$, we employ a Girsanov change of measure (controlling the KL-divergence), given as Lemma 5.2, followed by Corollary 4.5 (to relate back to 2-Wasserstein distance), as in [22]. This procedure relies crucially on the exponential integrability of the diffusion (2.12), which we prove as Lemma 7.3. To handle lack of measure absolute continuity, as discussed below, we must introduce an intermediate process (with law $\gamma_{k\epsilon}$), perform the above procedure on the error between $\gamma_{k\epsilon}$ and $\nu_{k\epsilon}$, then bound $\mathcal{W}_2(\pi_k, \nu_{k\epsilon}) \leq \mathcal{W}_2(\pi_k, \gamma_{k\epsilon}) + \mathcal{W}_2(\gamma_{k\epsilon}, \nu_{k\epsilon})$. The result providing a bound on $\mathcal{W}_2(\pi_k, \gamma_{k\epsilon})$, completing this approach, is given as Lemma 5.1.

To bound $W_2(\nu_{k\epsilon}, \pi_{\infty})$, we first show that π_{∞} satisfies a logarithmic-Sobolev inequality, by satisfying the conditions of Proposition 4.4 [5]. This result is given as Proposition 5.3. We then apply exponential decay of entropy [2], given as Lemma 4.1, and the Otto-Villani Theorem [1], given as 4.2. This procedure provides an exponentially decaying bound on $W_2(\nu_{k\epsilon}, \pi_{\infty})$.

5.1. 2-Wasserstein Bound for Diffusion Approximation. Here we obtain a bound on $W_2(\pi_k, \nu_{k\epsilon})$. Consider the continuous-time interpolation of the process (2.9):

(5.1)
$$\bar{\alpha}(t) = \alpha_0 - \int_0^t \left[K_\Delta(\theta_{\bar{s}}, \bar{\alpha}(\bar{s})) \frac{\beta}{2} \hat{\nabla} J(\theta_{\bar{s}}) + \nabla \pi_{0,\gamma}(\bar{\alpha}(\bar{s})) \right] \pi_{0,\gamma}(\bar{\alpha}(\bar{s})) ds + \int_0^t \pi_{0,\gamma}(\bar{\alpha}(\bar{s})) dW(s)$$

where $\bar{s} = \lfloor s/\epsilon \rfloor \epsilon$, and $\theta_{\bar{s}} := \theta_k$ for $k = \lfloor s/\epsilon \rfloor$. Note that, for each k, $\bar{\alpha}(k\epsilon)$ and α_k have the same probability law π_k . We aim to relate this process to the diffusion (2.12) through a Girsanov change of measure; but the process (5.1) is not Markovian and is therefore not an Itô diffusion. However, by the result of [11], the process $\bar{\alpha}(t)$ has the



Figure 2: Theorem 3.9 proof structure. First the 2-Wasserstein distance between discretetime algorithm (2.9) (with measure π_k) and continuous-time diffusion (2.12) (with measure $\nu_{k\epsilon}$) is bounded. We must introduce an intermediate process (with law $\gamma_{k\epsilon}$). Lemma 5.1 bounds the Wasserstein distance between π_k and $\gamma_{k\epsilon}$. Lemma 5.2 bounds the KL-divergence between π_k and $\gamma_{k\epsilon}$. Corollary 4.5 is then used, along with Lemma 7.3 to relate this KL bound to a 2-Wasserstein bound. Proposition 5.3 is the key tool in bounding $W_2(\nu_{k\epsilon}, \pi_{\infty})$, establishing that π_{∞} satisfies a log-Sobolev inequality. We then employ exponential decay of entropy (Lemma 4.1) and the Otto-Villani Theorem (Lemma 4.2) to obtain exponential decay of $W_2(\nu_{k\epsilon}, \pi_{\infty})$.

same one-time marginals as the Itô process Y(t), where

(5.2)
$$Y(t) = \alpha_0 - \int_0^t g_s(\theta_{\bar{s}}, Y(s)) ds + \int_0^t \mathbb{E} \left[\pi_{0,\gamma}(\bar{\alpha}(\bar{s})) \, | \, \bar{\alpha}(s) = Y(s) \right] dW(s)$$
$$g_s(\theta_{\bar{s}}, Y(s)) = \mathbb{E} \left[\left(K_\Delta(\theta_{\bar{s}}, \bar{\alpha}(\bar{s})) \frac{\beta}{2} \hat{\nabla} J(\theta_{\bar{s}}) + \nabla \pi_{0,\gamma}(\bar{\alpha}(\bar{s})) \right) \pi_{0,\gamma}(\bar{\alpha}(\bar{s})) \, \left| \, \bar{\alpha}(s) = Y(s) \right] \right]$$

i.e., $\text{Law}(Y(k\epsilon)) = \pi_k \forall k \in \mathbb{N}$, and it is apparent that (5.2) is Markovian. However, we cannot apply Girsanov's formula to relate (5.2) and (2.12) because the volatility functions are different; so the measures π_k and $\nu_{k\epsilon}$ are not absolutely continuous.

To solve this, we introduce the intermediate process

(5.3)
$$X(t) = \alpha_0 - \int_0^t \hat{g}_s(\theta_{\bar{s}}, X(s)) ds + \int_0^t \pi_{0,\gamma}(X(s)) dW(s)$$

where $\hat{g}_s(\theta_{\bar{s}}, X(s)) = \left(K_{\Delta}(\theta_{\bar{s}}, X(s))\frac{\beta}{2}\hat{\nabla}J(\theta_{\bar{s}}) + \nabla\pi_{0,\gamma}(X(s))\right)\pi_{0,\gamma}(X(s))$. Let $\gamma_{k\epsilon}$ denote the law of (5.3) at time $t = k\epsilon$. This process (5.3) is similar enough to (5.2) to allow a tractable bound on $\mathcal{W}_2(\pi_k, \gamma_{k\epsilon})$, and since (5.3) has the same volatility function as (2.12) we can relate these two via Girsanov's formula to obtain a desirable bound on $\mathcal{W}_2(\gamma_{k\epsilon}, \nu_{k\epsilon})$. Then we simply bound $\mathcal{W}_2(\pi_k, \nu_{k\epsilon}) \leq \mathcal{W}_2(\pi_k, \gamma_{k\epsilon}) + \mathcal{W}_2(\gamma_{k\epsilon}, \nu_{k\epsilon})$.

The following Lemma provides a bound on $\mathcal{W}_2(\pi_k, \gamma_{k\epsilon})$.

LEMMA 5.1. Fixing the step size ϵ and time horizon $k\epsilon$, take the kernel scale parameter Δ small enough to satisfy (3.4), and sampling distribution scale parameter γ small enough to satisfy (3.5). Then we have

$$\mathcal{W}_2(\pi_k, \gamma_{k\epsilon}) \le 6(k\epsilon)\epsilon\sqrt{12C_0 + 3 + 3\sqrt{2(k\epsilon)\epsilon}}$$

where C_0 is a constant provided in supplementary material 7.1, M_{θ} is a uniform bound on $\mathbb{E}\|\theta_k\|^2$, see Lemma 7.4, $L_{\nabla J}$ is the Lipchitz constant for ∇J , see Assumption 3.1, $B = \|\nabla J(0)\|$, and ζ is the uniform noise variance bound in Assumption 3.3. Now, the following Lemma provides a bound on $D(\gamma_{k\epsilon}\|\nu_{k\epsilon})$, the Kullback-Leibler (KL) divergence between measures $\gamma_{k\epsilon}$ and $\nu_{k\epsilon}$, using a similar Girsanov change of measure as presented in [22].

LEMMA 5.2. Fixing the step size ϵ and time horizon $k\epsilon$, taking the kernel scale parameter Δ small enough to satisfy (3.4), and sampling distribution scale parameter γ small enough to satisfy (3.5). Then we have:

$$D(\gamma_{k\epsilon} \| \nu_{k\epsilon}) \le (k\epsilon)^3 \epsilon^3 \left[4\beta L_{\nabla J}^2 \left(72C_0 + 6\sqrt{C_0} + 18 + \sqrt{2} \right) \right] + (k\epsilon) \epsilon \left(2L_J^2 + 4C_0 \right)$$

Now we relate this KL divergence to a Wasserstein distance through Corollary 4.5, with $\mu = \gamma_{k\epsilon}$, $\nu = \nu_{k\epsilon}$, $\lambda = 1$. By Lemma 7.3 we have exponential integrability of the diffusion 2.12:

$$\log \int_{\mathbb{R}^N} e^{\lambda \|w\|^2} \nu_{k\epsilon}(dw) \le \kappa_0^{\gamma} + ((\beta b + N)2\epsilon + 2I')k\epsilon$$

Now, since $k\epsilon \ge 1$ by (3.6) and (3.3), and $\kappa_0^{\gamma} \le \kappa_0$, we can bound

$$\mathcal{W}_2(\gamma_{k\epsilon},\nu_{k\epsilon}) \le 2\sqrt{\frac{3}{2} + (\kappa_0 + (\beta b + N)2\epsilon + 2I')k\epsilon} \left(\sqrt{D(\gamma_{k\epsilon}||\nu_{k\epsilon})} + (D(\gamma_{k\epsilon}||\nu_{k\epsilon}))^{1/4}\right)$$

Applying Lemma 5.2, we have:

$$\mathcal{W}_2(\gamma_{k\epsilon},\nu_{k\epsilon}) \le 4\sqrt{\frac{3}{2} + C_1 k\epsilon} \sqrt{k\epsilon} \sqrt{\epsilon} \left(4\sqrt{\beta L_{\nabla J}^2 C_2} + 2\sqrt{2L_J^2 + 4C_0} \right)$$

with C_0, C_1, C_2 in supplementary material 7.1. So finally,

(5.4)

$$\mathcal{W}_{2}(\pi_{k},\nu_{k\epsilon}) \leq k\epsilon\sqrt{\epsilon} \left[6\sqrt{12C_{0}+3} + 3\sqrt{2} + 4\sqrt{3/2 + C_{1}} \left(4\sqrt{C_{2}} + 2\sqrt{2L_{J}^{2} + 4C_{0}} \right) \right]$$

So we achieve a discretization error bound $\mathcal{W}_2(\pi_k, \nu_{k\epsilon})$ which scales as $\mathcal{O}(k\epsilon\sqrt{\epsilon})$. In fact, this is tighter than the bound obtained in [22], which scales as $\mathcal{O}(k\epsilon\epsilon^{1/4})$. We represent this bound in terms of distinct units $k\epsilon$ and $\sqrt{\epsilon}$ (rather than $k\epsilon^{3/2}$) since in our final analysis we will take $k\epsilon$ large enough (but fixed), then ϵ small enough, so that $\mathcal{W}_2(\pi_k, \nu_{k\epsilon})$ decreases arbitrarily. We will need to first take $k\epsilon$ large enough to control the diffusion (2.12) convergence to the Gibbs measure. The following presents this convergence in terms of exponentially decaying distance $\mathcal{W}_2(\nu_{k\epsilon}, \pi_{\infty})$.

5.2. 2-Wasserstein Distance for Diffusion Convergence. Here we describe the method to bound $W_2(\nu_{k\epsilon}, \pi_{\infty})$. The strategy is as follows:

- i) Show that π_{∞} satisfies a logarithmic-Sobolev inequality.
- ii) Apply exponential decay of entropy, given as Lemma 4.1, with the relative entropy bound in Lemma 7.2, to derive a bound on $D(\nu_{k\epsilon} || \pi_{\infty})$
- iii) Apply the Otto-Villani Theorem, given as Lemma 4.2, to relate this to a bound on $W_2(\nu_{k\epsilon}, \pi_{\infty})$.

We accomplish (i) in the following proposition, establishing that the Gibbs measure π_{∞} satisfies a log-Sobolev inequality:

PROPOSITION 5.3. For β satisfying Assumption 3.8, the Gibbs measure π_{∞} satisfies a logarithmic Sobolev inequality with constant c_{LS} :

(5.5)

$$0 \le c_{LS} \le \frac{2\beta L_{\nabla J}}{\gamma} + \frac{2}{\beta L_{\nabla J}} + \frac{1}{\lambda} \left(\frac{2\beta L_{\nabla J}}{\gamma} \left(\kappa + \gamma \left(\kappa_0 + \frac{(\beta b + N)\bar{\pi}_{0,\gamma} + 2I}{(m\beta)\bar{\pi}_{0,\gamma}} \right) \right) + 2 \right)$$

where

$$\begin{aligned} \frac{1}{\lambda} &\leq \frac{1}{2\kappa} \left(1 + \frac{4C\kappa^2}{\gamma} \exp\left(\beta \left(\frac{(L_{\nabla J} + B)\kappa}{\gamma} + A + B\right)\right) \right) \\ \kappa &= \left(\frac{1}{2}\beta mN + \beta mI\right) + \frac{1}{2} \left[\beta^2 mb + (\beta mM)^2\right], \quad \gamma = \frac{1}{2} \left((\beta m)^2 + \left(1 - \frac{1}{\bar{\pi}_0^2 + 1}\right)\right) \end{aligned}$$

and $\bar{\pi}_{0,\gamma} = \sup_x \pi_{0,\gamma}(x), \ \bar{\pi}_0 = \sup_x \pi_0(x).$

Proof Sketch: The full proof is available in supplementary material 7.4. The key tool we use is the main Theorem in [5], reproduced as Proposition 4.4. To satisfy condition (1) of Proposition 4.4 we show that the Lyapunov function

$$V(w) = \exp\left(\frac{\beta m \|w\|^2}{2(\bar{\pi}^2_{0,\gamma} + 1)}\right)$$

and the infinitesimal generator (4.3) satisfy (4.10), with κ and γ given in (5.6). Then, Proposition 4.3 is used to show that condition (2) is satisfied. Condition (3) is satisfied with $K = \beta L_{\nabla J}$ by assumption 3.1.

Now since $D(\nu_0||\pi_\infty) = D(\pi_0||\pi_\infty) < \infty$ by Lemma 7.2, we can apply the exponential decay of entropy (Lemma 4.1) to obtain

(5.7)
$$D(\nu_t || \pi_{\infty}) \le D(\pi_{0,\gamma} || \pi_{\infty}) e^{-2t/\beta c_{LS}}$$

Then by the Otto-Villani Theorem and Lemma 7.2, we have

(5.8)
$$\mathcal{W}_2(\nu_t, \pi_\infty) \le \sqrt{2c_{LS}\bar{D}_0^{\gamma}} e^{-t/\beta c_{LS}}$$

where \bar{D}_0^{γ} is the relative entropy bound given in (7.2) and c_{LS} is bounded in (5.5).

5.3. Controlling the 2-Wasserstein Distance. Combining the bounds (5.4) and (5.8) yields

$$\mathcal{W}_{2}(\pi_{k},\pi_{\infty}) \leq k\epsilon\sqrt{\epsilon} \left[6\sqrt{12C_{0}+3} + 3\sqrt{2} + 4\sqrt{\frac{3}{2}+C_{1}} \left(4\sqrt{C_{2}} + 2\sqrt{2L_{J}^{2}+4C_{0}} \right) \right]$$

$$(5.9) + \sqrt{2c_{LS}\bar{D}_{0}^{\gamma}}e^{-k\epsilon/\beta c_{LS}}$$

The strategy to control (5.9) is to take $k\epsilon$ large enough so that the exponential term dies away, then (fixing $k\epsilon$) take ϵ small enough so that the first term decreases arbitrarily. However, we encounter a subtle problem: the term \bar{D}_0^{γ} may depend inconveniently on γ , and thus on ϵ , since we take γ satisfying (3.5) in order to obtain Lemmas 5.1, 5.2.

Let us investigate this. Lemma 7.2, with $\gamma \leq 1$ and $\bar{\pi}_{0,\gamma}$ expanded, gives

$$\bar{D}_0^{\gamma} \le \log(\bar{\pi}_0) + \log\frac{1}{\gamma^N} + \frac{N}{2}\log\frac{3\pi}{m\beta} + \frac{\beta b}{2}\log3 + \beta\left(\frac{L_{\nabla J}}{3}\kappa_0 + B\sqrt{\kappa_0} + A\right)$$
18

so we see that \bar{D}_0^{γ} depends on γ as $N \log \left(\frac{1}{\gamma}\right)$. Observe that taking $k\epsilon$ as (3.3) and ϵ as (3.2) yields

$$\mathcal{W}_{2}(\pi_{k},\pi_{\infty}) \leq \delta \left[6\sqrt{12C_{0}+3} + 3\sqrt{2} + 4\sqrt{\frac{3}{2}+C_{1}} \left(4\sqrt{C_{2}} + 2\sqrt{2L_{J}^{2}+4C_{0}} \right) \right]$$

$$(5.10) + \delta \sqrt{2c_{LS}N \log\left(\frac{1}{\gamma}\right)} + \delta \sqrt{2c_{LS}C_{3}}$$

where $C_3 := \log(\bar{\pi}) + \frac{N}{2}\log\frac{3\pi}{m\beta} + \frac{\beta b}{2}\log 3 + \beta \left(\frac{L_{\nabla J}}{3}\kappa_0 + B\sqrt{\kappa_0} + A\right)$. Then, since $\gamma \in [\epsilon^2, \epsilon^{3/2}]$ and $\epsilon \leq \left(\frac{\delta}{\log(1/\delta)}\right)^2$ we have

$$\log(\frac{1}{\gamma}) \le \log\left(\frac{1}{\epsilon^2}\right) \le \log\left(\frac{\log(1/\delta)}{\delta}\right)^4 \le \log\left(\frac{1}{\delta^5}\right) = 5\log\left(1/\delta\right)$$

where we use that $\left(\frac{\log(1/\delta)}{\delta}\right)^4 \leq \delta^{-5}$ for all $\delta \leq 1$, satisfied by the feasible δ range (3.6). Then we obtain the bound displayed in Theorem 3.9.

6. Conclusion. We derived non-asymptotic (finite-sample) bounds for a passive stochastic gradient Langevin dynamics algorithm. These results complement recent asymptotic weak convergence analysis of the passive Langevin algorithm in [16]. The passive Langevin algorithm analyzed in this paper uses sequential evaluations of a stochastic gradient descent by an external agent (forward learner), and reconstructs the cost function being optimized. Thus it achieves *real-time* (adaptive) inverse reinforcement learning, in that we (the inverse learner) reconstruct the cost function while it is in the process of being optimized. Specifically, we have provided finite-sample bounds on the 2-Wasserstein distance between the sample distribution induced by our algorithm and the Gibbs measure encoding the cost function to be reconstructed. Our paper builds on the seminal paper [22] and utilizes techniques in the analysis of Markov Diffusion Operators [2] to achieve the bound.

REFERENCES

- D. BAKRY, F. BARTHE, P. CATTIAUX, AND A. GUILLIN, A simple proof of the Poincaré inequality for a large class of probability measures, (2008).
- [2] D. BAKRY, I. GENTIL, M. LEDOUX, ET AL., Analysis and geometry of Markov diffusion operators, vol. 103, Springer, 2014.
- [3] F. BOLLEY AND C. VILLANI, Weighted Csiszár-Kullback-Pinsker inequalities and applications to transportation inequalities, in Annales de la Faculté des sciences de Toulouse: Mathématiques, vol. 14, 2005, pp. 331–352.
- [4] D. BROWN, W. GOO, P. NAGARAJAN, AND S. NIEKUM, Extrapolating beyond suboptimal demonstrations via inverse reinforcement learning from observations, in International conference on machine learning, PMLR, 2019, pp. 783–792.
- [5] P. CATTIAUX, A. GUILLIN, AND L. WU, A note on Talagrand's transportation inequality and logarithmic sobolev inequality, arXiv preprint arXiv:0810.5435, (2008).
- [6] X. CHENG, N. S. CHATTERJI, P. L. BARTLETT, AND M. I. JORDAN, Underdamped Langevin MCMC: A non-asymptotic analysis, in Conference on learning theory, PMLR, 2018, pp. 300–323.
- [7] T.-S. CHIANG, C.-R. HWANG, AND S. J. SHEU, Diffusion for global optimization in rⁿ, SIAM Journal on Control and Optimization, 25 (1987), pp. 737–753.
- [8] A. S. DALALYAN AND A. KARAGULYAN, User-friendly guarantees for the Langevin Monte Carlo with inaccurate gradient, Stochastic Processes and their Applications, 129 (2019), pp. 5278– 5311.

- H. DJELLOUT, A. GUILLIN, AND L. WU, Transportation cost-information inequalities and applications to random dynamical systems and diffusions, The Annals of Probability, 32 (2004), pp. 2702–2732.
- [10] S. B. GELFAND AND S. K. MITTER, Simulated annealing type algorithms for multivariate optimization, Algorithmica, 6 (1991), pp. 419–436.
- [11] I. GYÖNGY, Mimicking the one-dimensional marginal distributions of processes having an Itô differential, Probability theory and related fields, 71 (1986), pp. 501–516.
- [12] D. HADFIELD-MENELL, S. RUSSELL, P. ABBEEL, AND A. DRAGAN, Cooperative inverse reinforcement learning, in Advances in neural information processing systems, IEEE, 2016, pp. 3909–3917.
- [13] W. HARDLE AND R. NIXDORF, Nonparametric sequential estimation of zeros and extrema of regression functions, IEEE transactions on information theory, 33 (1987), pp. 367–372.
- [14] P. KAIROUZ, H. B. MCMAHAN, B. AVENT, A. BELLET, M. BENNIS, A. N. BHAGOJI, K. BONAWITZ, Z. CHARLES, G. CORMODE, R. CUMMINGS, ET AL., Advances and open problems in federated learning, Foundations and Trends[®] in Machine Learning, 14 (2021), pp. 1–210.
- [15] I. KARATZAS AND S. E. SHREVE, Brownian motion and stochastic calculus, vol. 113, Springer Science & Business Media, 1991.
- [16] V. KRISHNAMURTHY AND G. YIN, Langevin dynamics for adaptive inverse reinforcement learning of stochastic gradient algorithms., J. Mach. Learn. Res., 22 (2021), pp. 121–1.
- [17] V. KRISHNAMURTHY AND G. YIN, Multikernel passive stochastic gradient algorithms and transfer learning, IEEE Transactions on Automatic Control, 67 (2022), pp. 1792–1805.
- [18] A. KROGH AND J. HERTZ, A simple weight decay can improve generalization, Advances in neural information processing systems, 4 (1991).
- [19] H. J. KUSHNER AND G. YIN, Stochastic Approximation Algorithms and Recursive Algorithms and Applications, Springer-Verlag, 2nd ed., 2003.
- [20] A. V. NAZIN, B. T. POLYAK, AND A. B. TSYBAKOV, Passive stochastic approximation, Avtomatika i Telemekhanika, (1989), pp. 127–134.
- [21] A. Y. NG, S. RUSSELL, ET AL., Algorithms for inverse reinforcement learning., in Icml, vol. 1, 2000, p. 2.
- [22] M. RAGINSKY, A. RAKHLIN, AND M. TELGARSKY, Non-convex learning via stochastic gradient Langevin dynamics: a nonasymptotic analysis, in Conference on Learning Theory, PMLR, 2017, pp. 1674–1703.
- [23] P. RÉVÉSZ, How to apply the method of stochastic approximation in the non-parametric estimation of a regression function, Statistics: A Journal of Theoretical and Applied Statistics, 8 (1977), pp. 119–126.
- [24] D. E. RUMELHART, G. E. HINTON, AND R. J. WILLIAMS, Learning representations by backpropagating errors, nature, 323 (1986), pp. 533–536.
- [25] O. STRAMER AND R. TWEEDIE, Langevin-type models i: Diffusions with given stationary distributions and their discretizations, Methodology and Computing in Applied Probability, 1 (1999), pp. 283–306.
- [26] R. S. SUTTON, D. MCALLESTER, S. SINGH, AND Y. MANSOUR, Policy gradient methods for reinforcement learning with function approximation, Advances in neural information processing systems, 12 (1999).
- [27] C. VILLANI, Topics in optimal transportation, vol. 58, American Mathematical Soc., 2021.
- [28] M. WELLING AND Y. W. TEH, Bayesian learning via stochastic gradient Langevin dynamics, in Proceedings of the 28th international conference on machine learning (ICML-11), 2011, pp. 681–688.
- [29] P. XU, J. CHEN, D. ZOU, AND Q. GU, Global convergence of Langevin dynamics based algorithms for nonconvex optimization, Advances in Neural Information Processing Systems, 31 (2018).
- [30] G. YIN AND K. YIN, Passive stochastic approximation with constant step size and window width, IEEE transactions on automatic control, 41 (1996), pp. 90–106.
- [31] G. G. YIN AND K. YIN, Passive stochastic approximation with constant step size and window width, IEEE transactions on automatic control, 41 (1996), pp. 90–106.
- [32] B. D. ZIEBART, A. L. MAAS, J. A. BAGNELL, A. K. DEY, ET AL., Maximum entropy inverse reinforcement learning., in Aaai, vol. 8, Chicago, IL, USA, 2008, pp. 1433–1438.

7. Supplementary Material.

Glossary of Symbols.

SGD (2.8) iterate, $\theta_k \in \mathbb{R}^N$ θ_k PSGLD (2.9) iterate, $\alpha_k \in \mathbb{R}^N$ α_k ϵ PSGLD step size sampling distribution (2.7), maximum value: $\bar{\pi}_{0,\gamma}$ $\pi_{0,\gamma}$ PSGLD inverse temperature parameter β KPSGLD kernel function Δ PSGLD kernel function scale parameter cost function, $J : \mathbb{R}^N \to \mathbb{R}_+$ J Gibbs measure $(\pi_{\infty} \propto \exp(-\beta J))$ π_{∞} solution of Itô diffusion (2.12) $\alpha(t)$ W(t)standard Brownian motion diffusion (2.12) log-Sobolev constant c_{LS} diffusion (2.12) Poincarê constant c_P ∇J Lipschitz constant $L_{\nabla J}$ (m, b)J dissipativity constants uniform stochastic gradient $\hat{\nabla} J(\cdot)$ variance bound Č L_{π_0} π_0 Lipschitz constant uniform SGD bound, $M_{\theta} := \sup_{k \ge 0} \mathbb{E} \|\theta_k\|^2$ $\kappa_0^{\gamma} := \log \mathbb{E}_{\pi_{0,\gamma}} \left[\exp(\|x\|^2) \right] < \infty \ (\kappa_0 := \kappa_0^{\gamma}|_{\gamma=1})$ M_{θ} κ_0^{γ} (7.6) inner product bound 1 Ι I'(7.6) inner product bound 2 A|J(0)|B $\|\nabla J(0)\|$ $\|\cdot\|$ l_2 norm

7.1. Bound Constants.

$$C_{0} := 3L_{\nabla J}^{2}(M_{\theta} + 2B^{2}M_{\theta}) + B^{2} + \zeta$$

$$C_{1} := \kappa_{0} + (\beta b + N)2\epsilon + 2I')$$

$$C_{2} := \beta L_{\nabla J}^{2} \left(72C_{0} + 6\sqrt{C_{0}} + 18 + \sqrt{2}\right)$$

$$C_{3} := \log(\bar{\pi}) + \frac{N}{2}\log\frac{3\pi}{m\beta} + \frac{\beta b}{2}\log3 + \beta\left(\frac{L_{\nabla J}}{3}\kappa_{0} + B\sqrt{\kappa_{0}} + A\right)$$

$$C_{4} := \left[6\sqrt{12C_{0} + 3} + 3\sqrt{2} + 4\left(\frac{3}{2} + C_{1}\right)^{1/2}\left(4\sqrt{C_{2}} + 2\sqrt{2L_{J}^{2} + 4C_{0}}\right)\right]$$

$$M_{\theta} = \kappa_{0} + 2\left(1 \vee \frac{1}{m}\right)\left(b + 2B^{2}\right)$$

7.2. Technical Results. Here we present several technical Lemmas which are necessary for the results derived in Section 5.

The proofs for all of these can be found in supplementary material 7.3. We denote $\bar{\pi}_0 := \sup_x \pi_0(x)$ and $\bar{\pi}_{0,\gamma} := \sup_x \pi_{0,\gamma}(x)$. $A = ||J(0)||, B = ||\nabla J(0)||$, and I, I' are constants provided in Lemma 7.11.

LEMMA 7.1 ($\pi_{0,\gamma}$ exponential integrability). For all $\gamma \leq 1$, $\pi_{0,\gamma}$ has a bounded 21 and strictly positive density with respect to the Lebesgue measure on \mathbb{R}^N , and

(7.1)
$$\kappa_0^{\gamma} := \log \int_{\mathbb{R}^N} e^{\|x\|^2} d\pi_{0,\gamma}(x) < \infty$$

and denote $\kappa_0 := \kappa_0^{\gamma}|_{\gamma=1}$ so that $\kappa_0^{\gamma} \leq \kappa_0 \ \forall \gamma \leq 1$.

LEMMA 7.2 (relative entropy bound).

(7.2)
$$\bar{D}_{0}^{\gamma} := D(\pi_{0,\gamma} || \pi_{\infty}) \leq \log \bar{\pi}_{0,\gamma} + \frac{N}{2} \log \frac{3\pi}{m\beta} + \frac{\beta b}{2} \log 3 + \beta \left(\frac{L_{\nabla J}}{3} \kappa_{0}^{\gamma} + B \sqrt{\kappa_{0}^{\gamma}} + A \right)$$

LEMMA 7.3 (exponential integrability of Langevin diffusion).

$$\log \mathbb{E}[e^{\|\alpha(t)\|^2}] \le \kappa_0^{\gamma} + ((\beta b + N)2\epsilon + 2I')t$$

where κ_0 is given in (7.1) and I' is given in (7.6).

LEMMA 7.4 (uniform L^2 bound on SGD). For $\eta \in (0, 1 \land \frac{m}{4L_{\nabla J}^2})$,

(7.3)
$$\sup_{k\geq 0} \mathbb{E}\|\theta_k\|^2 \leq \kappa_0 + 2\left(1 \vee \frac{1}{m}\right) \left(b + 2B^2\right) =: M_\theta$$

LEMMA 7.5 (L^2 bound on Langevin diffusion).

$$\mathbb{E}\|\alpha(t)\|^2 \le \kappa_0^{\gamma} + \frac{(\beta b + N)\bar{\pi}_{0,\gamma} + 2I}{(m\beta)\bar{\pi}_{0,\gamma}}$$

LEMMA 7.6. Taking

(7.4)
$$\Delta \leq \inf_{x \in [\epsilon, \hat{K}_{\epsilon}]} \frac{K^{-1}(\frac{\hat{K}_{1}\sqrt{2\pi}}{2\epsilon}e^{x^{2}/2})}{K^{-2}(x\epsilon^{2N})}$$

gives

$$\mathbb{E} \|K_{\Delta}(\theta_k, \alpha_k) \hat{\nabla} J(\theta_k)\|^2 \le 12\epsilon (L_{\nabla J}^2(M_{\theta} + 2B^2 M_{\theta}) + B^2 + \zeta)$$

LEMMA 7.7. Taking π_0 such that $\bar{\pi}_0 = 1$, and

(7.5)
$$\gamma \le \epsilon^{3/2}$$

gives

$$\mathbb{E} \|\pi_{0,\gamma}(\alpha_k)\|^2 \le \epsilon, \quad \mathbb{E} \|\nabla \pi_{0,\gamma}(\alpha_k)\|^2 \le \epsilon$$

LEMMA 7.8. Let R be an N-dimensional random variable on the same probability space as $\bar{\alpha}(s)$. Then

$$\mathbb{E}\|\pi_{0,\gamma}(\bar{\alpha}(s))R\|^2 \le 2\mathbb{E}\left[|\pi_{0,\gamma}(\bar{\alpha}(s))|^2\right]\mathbb{E}\left[\|R\|^2\right]$$

LEMMA 7.9 ($\pi_{0,\gamma}$ Quadratic Decay). $\pi_{0,\gamma}$ has tail value decay $\mathcal{O}(||x||^{-2})$ and differential decay $\mathcal{O}(||x||^{-1})$. Specifically,

$$\exists M \in \mathbb{N}, \, \tilde{C} > 0 : \, \pi_{0,\gamma}(x) \le \frac{2}{(\beta m^*)^2 \|x\|^2}, \, \|\nabla \pi_{0,\gamma}(x)\| \le \frac{\tilde{C}}{\|x\|} \, \forall \|x\| > M$$

LEMMA 7.10 (quadratic bounds on J). For all $w \in \mathbb{R}^d$,

$$\|\nabla J(w)\| \le L_{\nabla J} \|w\| + B$$

and

$$\frac{m}{3} \|w\|^2 - \frac{b}{2} \le J(w) \le \frac{L_{\nabla J}}{2} \|w\|^2 + B\|w\| + A$$

-

LEMMA 7.11 (Uniform Gradient Inner Product Bound). We have

(7.6)

$$\begin{aligned} \exists I > 0 : \|\langle x, \nabla \pi_{0,\gamma}(x) \rangle\| &\leq I \ \forall x \in \mathbb{R}^N, \gamma \in \mathbb{R}_+ \\ \exists I' > 0 : \|\langle x, \pi_{0,\gamma}(x) \nabla \pi_{0,\gamma}(x) \rangle\| &\leq I' \ \forall x \in \mathbb{R}^N, \gamma \in \mathbb{R}_+ \end{aligned}$$

7.3. Proofs of Technical Lemmas.

Proof of Lemma 7.1.

Proof. This follows from Assumption 3.4.

Proof of Lemma 7.2.

Proof. Recall $\pi_{\infty}(w) := \frac{1}{\Lambda} \exp(-\beta J(w))$, where $\Lambda = \int_{\mathbb{R}^N} \exp(-\beta J(w)) dw$. Since $\pi_{\infty} > 0$ everywhere, we can write

$$D(\pi_{0,\gamma}||\pi_{\infty}) = \int_{\mathbb{R}^{N}} \pi_{0,\gamma}(x) \log\left(\frac{\pi_{0,\gamma}(x)}{\pi_{\infty}(x)}\right) dx$$

$$= \int_{\mathbb{R}^{N}} \pi_{0,\gamma}(x) \log \pi_{0,\gamma}(x) dx + \log \Lambda + \beta \int_{\mathbb{R}^{N}} \pi_{0,\gamma}(x) J(x) dx$$

$$\leq \log \|\pi_{0,\gamma}\|_{\infty} + \log \Lambda + \beta \int_{\mathbb{R}^{N}} \pi_{0,\gamma}(x) J(x) dx$$

$$\leq \log \bar{\pi}_{0,\gamma} + \log \Lambda + \beta \int_{\mathbb{R}^{N}} \pi_{0,\gamma}(x) J(x) dx$$

First let us upper bound the normalization constant:

$$\begin{split} \Lambda &= \int_{\mathbb{R}^N} e^{-\beta J(x)} dx \\ &\leq e^{\frac{1}{2}\beta b \log 3} \int_{\mathbb{R}^N} e^{-\frac{m\beta \|x\|^2}{3}} dx \\ &= 3^{\beta b/2} \left(\frac{3\pi}{m\beta}\right)^{N/2} \end{split}$$

where the inequality follows from Lemma 7.10. Thus,

$$\log\Lambda \leq \frac{N}{2}\log\frac{3\pi}{m\beta} + \frac{\beta b}{2}\log 3$$

By Lemma 7.10 we also have

$$\int_{\mathbb{R}^N} J(x)\pi_{0,\gamma}(x)dx \leq \int_{\mathbb{R}^N} \pi_{0,\gamma}dx \left(\frac{L_{\nabla J}}{3} \|x\|^2 + B\|x\| + A\right)$$
$$\leq \frac{L_{\nabla J}}{3}\kappa_0^{\gamma} + B\sqrt{\kappa_0^{\gamma}} + A$$

Thus

$$D(\pi_{0,\gamma}||\pi_{\infty}) \le \log \bar{\pi}_{0,\gamma} + \frac{N}{2}\log \frac{3\pi}{m\beta} + \frac{\beta b}{2}\log 3 + \beta \left(\frac{L_{\nabla J}}{3}\kappa_0^{\gamma} + B\sqrt{\kappa_0^{\gamma}} + A\right) \quad \Box$$

Proof of Lemma 7.3.

Proof. Let us denote $\alpha_t := \alpha(t)$ for notational convenience, and define $L(t) := e^{\|\alpha_t\|^2}$. Similarly denote $L_t = L(t)$. By Itô's Lemma we have

$$dL_t = \{ (\nabla_{\alpha_t} L_t)^T \boldsymbol{\mu}_t + \frac{1}{2} \operatorname{Tr}[\boldsymbol{G}_t^T (\boldsymbol{H}_{\alpha_t} L_t) \boldsymbol{G}_t] \} dt + (\nabla_{\alpha_t} L_t)^T \boldsymbol{G}_t dW(t) \}$$

where from (2.12) we have

$$\boldsymbol{\mu}_t = -\frac{\beta}{2} \pi_{0,\gamma}^2(\alpha_t) \nabla J(\alpha_t) - \pi_{0,\gamma}(\alpha_t) \nabla \pi_{0,\gamma}(\alpha_t), \ \boldsymbol{G}_t = \pi_{0,\gamma}(\alpha_t)$$

Thus,

$$(\nabla_{\alpha_t} L_t)^T \boldsymbol{\mu}_t = -\beta \langle \alpha_t L_t, \pi_{0,\gamma}^2(\alpha_t) \nabla J(\alpha_t) \rangle - 2 \langle \alpha_t L_t, \pi_{0,\gamma}(\alpha_t) \nabla \pi_{0,\gamma}(\alpha_t) \rangle$$

and

$$\frac{1}{2} \text{Tr}[\boldsymbol{G}_t^T(H_{\alpha_t} L_t) \boldsymbol{G}_t] = \frac{1}{2} \text{Tr}[\pi_{0,\gamma}^2(\alpha_t) H_{\alpha_t} L_t] = \pi_{0,\gamma}^2(\alpha_t) (\|\alpha_t\|^2 L_t + NL_t)$$

and

$$(\nabla_{\alpha_t} L_t)^T \boldsymbol{G}_t = 2\alpha_t^* \pi_{0,\gamma}(\alpha_t) L_t$$

Putting these together and integrating,

$$\begin{split} L_t &= L(0) - \beta \int_0^t \langle \alpha_s L_s, \pi_{0,\gamma}^2(\alpha_s) \nabla J(\alpha_s) \rangle ds - 2 \int_0^t \langle \alpha_s L_s, \pi_{0,\gamma}(\alpha_s) \nabla \pi_{0,\gamma}(\alpha_s) \rangle ds \\ &+ \int_0^t \pi_{0,\gamma}^2(\alpha_s) (\|\alpha_s\|^2 L_s + NL_s) ds + \int_0^t 2\alpha_s^* \pi_{0,\gamma}(\alpha_s) L_s dW(s) \\ &= L(0) + \int_0^t (\pi_{0,\gamma}^2(\alpha_s) \|\alpha_s\|^2 - \beta \langle \alpha_s, \pi_{0,\gamma}^2(\alpha_s) \nabla J(\alpha_s) \rangle) L_s ds \\ &- 2 \int_0^t \langle \alpha_s L_s, \pi_{0,\gamma}(\alpha_s) \nabla \pi_{0,\gamma}(\alpha_s) \rangle ds + \int_0^t N L_s \pi_{0,\gamma}^2(\alpha_s) ds \\ &+ \int_0^t 2\alpha_s^* \pi_{0,\gamma}(\alpha_s) L_s dW(s) \end{split}$$

Now, from the dissipativity condition 3.2, we can obtain the following bound:

$$\begin{aligned} \pi_{0,\gamma}^2(\alpha_s)(\|\alpha_s\|^2 - \beta\langle\alpha_s, \nabla J(\alpha_s)\rangle) &\leq \pi_{0,\gamma}^2(\alpha_s)(\|\alpha_s\|^2 + \beta[-m\|\alpha_s\|^2 + b]) \\ &= \|\alpha_s\|^2(\pi_{0,\gamma}^2(\alpha_s) - \beta m) + \pi_{0,\gamma}^2(\alpha_s)\beta b \leq \pi_{0,\gamma}^2(\alpha_s)\beta b \end{aligned}$$

Making this substitution, we now work with

$$L_t \leq L(0) + (\beta b + N) \int_0^t \pi_{0,\gamma}^2(\alpha(s)) L_s ds - 2 \int_0^t \langle \alpha_s L_s, \pi_{0,\gamma}(\alpha_s) \nabla \pi_{0,\gamma}(\alpha_s) \rangle ds + \int_0^t 2\alpha_s^* \pi_{0,\gamma}(\alpha_s) L_s dW(s)$$

It can be shown (e.g., proof of Corollary 4.1 in [9]) that $\int_0^T \mathbb{E} \|L_t \alpha(t)\|^2 dt < \infty \ \forall T \ge 0$. Therefore the Itô integral $\int L_s \alpha_s^* dW(s)$ is a zero-mean martingale. Thus, taking expectations leaves us with

$$\begin{split} \mathbb{E}[L_t] &\leq \mathbb{E}[L(0)] + \mathbb{E}[(\beta b + N) \int_0^t \pi_{0,\gamma}^2(\alpha(s)) L_s ds] \\ &- \mathbb{E}[2 \int_0^t \langle \alpha_s, \pi_{0,\gamma}(\alpha_s) \nabla \pi_{0,\gamma}(\alpha_s) \rangle L_s ds] \\ &= \mathbb{E}[L(0)] + (\beta b + N) \int_0^t \mathbb{E}[\pi_{0,\gamma}^2(\alpha(s)) L_s] ds \\ &- 2 \int_0^t \mathbb{E} \langle \alpha_s, \pi_{0,\gamma}(\alpha_s) \nabla \pi_{0,\gamma}(\alpha_s) \rangle L_s ds] \\ &\leq \mathbb{E}[L(0)] + (\beta b + N) 2\epsilon \int_0^t \mathbb{E}[L_s] ds + 2I' \int_0^t \mathbb{E}[L_s] ds \\ &= \mathbb{E}[L(0)] + ((\beta b + N) 2\epsilon + 2I') \int_0^t \mathbb{E}[L_s] ds \\ &= e^{\kappa_0^{\gamma}} + ((\beta b + N) 2\epsilon + 2I') \int_0^t \mathbb{E}[L_s] ds \end{split}$$

By application of the Gronwall Inequality, we obtain

$$\mathbb{E}[L_t] \le \exp(\kappa_0^{\gamma} + \int_0^t ((\beta b + N)2\epsilon + 2I')ds$$
$$= \exp(\kappa_0^{\gamma} + ((\beta b + N)2\epsilon + 2I')t)$$

Thus,

$$\log \mathbb{E}[e^{\|\alpha(t)\|^2}] \le \kappa_0^\gamma + ((\beta b + N)2\epsilon + 2I')t$$

Proof of Lemma 7.4.

Proof. See Lemma 3 of [22], with $\beta = \infty$.

Proof of Lemma 7.5.

Proof. We consider the diffusion given by (2.12). Letting $Y(t) = \|\alpha(t)\|^2$, Itô's Lemma gives

$$dY(t) = \left[-2\langle \alpha(t), \frac{\beta}{2} \pi_{0,\gamma}^2(\alpha(t)) \nabla J(\alpha(t)) + \pi_{0,\gamma}(\alpha(t)) \nabla \pi_{0,\gamma}(\alpha(t)) \rangle + N \pi_{0,\gamma}^2(\alpha(t)) \right] dt + \pi_{0,\gamma}(\alpha(t)) \alpha(t)^* dW(t)$$

where $\alpha(t)^* dW(t) := \sum_{i=1}^N \alpha_i(t) dW_i(t)$. Letting $m := \frac{m\beta}{2} \bar{\pi}_{0,\gamma}^2$, we then form $d(e^{2mt} V(t)) = 2me^{2mt} V(t) + e^{2mt} dV(t)$

$$\begin{aligned} d(e^{2mt}Y(t)) &= 2me^{2mt}Y(t) + e^{2mt}dY(t) \\ &= \left[-2e^{2mt}\langle \alpha(t), \frac{\beta}{2}\pi_{0,\gamma}^2(\alpha(t))\nabla J(\alpha(t)) + \pi_{0,\gamma}(\alpha(t))\nabla \pi_{0,\gamma}(\alpha(t))\rangle \right. \\ &+ N\pi_{0,\gamma}^2(\alpha(t))e^{2mt} + 2me^{2mt}Y(t) \right] dt \\ &+ e^{2mt}\pi_{0,\gamma}(\alpha(t))\alpha(t)^* dW(t) \end{aligned}$$

Then integrating yields

$$Y(t) = e^{-2mt}Y(0) - 2\int_0^t e^{2m(s-t)} \langle \alpha(s), \frac{\beta}{2}\pi_{0,\gamma}^2(\alpha(s))\nabla J(\alpha(s)) \rangle_{25}$$

$$+ \pi_{0,\gamma}(\alpha(s))\nabla\pi_{0,\gamma}(\alpha(s))\rangle ds$$

$$+ 2m \int_0^t e^{2m(s-t)}Y(s)ds + \int_0^t N\pi_{0,\gamma}^2(\alpha(s))e^{2m(s-t)}ds$$

$$+ 2\int_0^t e^{2m(s-t)}\pi_{0,\gamma}(\alpha(s))\alpha(s)^*dW(s)$$

Then using the dissipativity condition 3.2 we get

$$\begin{split} Y(t) &\leq e^{-2mt}Y(0) + \beta \int_0^t \pi_{0,\gamma}^2(\alpha(s))e^{2m(s-t)}(b-mY(s))ds \\ &\quad + 2\int_0^t \pi_{0,\gamma}(\alpha(s))e^{2m(s-t)}Ids \\ &\quad + m\beta \int_0^t \bar{\pi}_{0,\gamma}^2 e^{2m(s-t)}Y(s)ds + \int_0^t N\bar{\pi}_{0,\gamma}^2 e^{2m(s-t)}ds \\ &\quad + 2\int_0^t e^{2m(s-t)}\pi_{0,\gamma}(\alpha(s))\alpha(s)^*dW(s) \\ &\leq e^{-2mt}Y(0) + \beta b\int_0^t \bar{\pi}_{0,\gamma}^2 e^{2m(s-t)}ds + 2I\int_0^t \bar{\pi}_{0,\gamma} e^{2m(s-t)}ds \\ &\quad + \int_0^t N\bar{\pi}_{0,\gamma}^2 e^{2m(s-t)}ds \end{split}$$

Then grouping terms and evaluating the integral yields

$$Y(t) \le e^{-2mt} Y(0) + \frac{(\beta b + N)\bar{\pi}_{0,\gamma}^2 + 2I\bar{\pi}_{0,\gamma}}{2m} \left(1 - e^{-2mt}\right) \\ + \bar{\pi}_{0,\gamma} \int_0^t e^{2m(s-t)} \alpha(s)^* dW(s)$$

Now taking expectations, and by the Martingale property of the Itô integral, we have

$$\mathbb{E}[\|\alpha(t)\|^{2}] \leq e^{-2mt} \mathbb{E}\|\alpha(0)\|^{2} + \frac{(\beta b + N)\bar{\pi}_{0,\gamma}^{2} + 2I\bar{\pi}_{0,\gamma}}{2m} \left(1 - e^{-2mt}\right)$$
$$\leq e^{-2mt} \mathbb{E}\|\alpha(0)\|^{2} + \frac{(\beta b + N)\bar{\pi}_{0,\gamma}^{2} + 2I\bar{\pi}_{0,\gamma}}{2m}$$

and from (7.1), and using $m = \frac{m\beta}{2} \bar{\pi}^2_{0,\gamma}$, and taking the maximum over t > 0 gives

$$\mathbb{E}\|\alpha(t)\|^2 \le \kappa_0^{\gamma} + \frac{(\beta b + N)\bar{\pi}_{0,\gamma} + 2I}{(m\beta)\bar{\pi}_{0,\gamma}} \qquad \Box$$

Proof of Lemma 7.6.

Proof.

$$\mathbb{E} \|K_{\Delta}(\theta_{k},\alpha_{k})\hat{\nabla}J(\theta_{k})\|^{2} = \mathbb{E} |K_{\Delta}(\theta_{k},\alpha_{k})|^{2} \|\hat{\nabla}J(\theta_{k})\|^{2}$$

$$\leq \mathbb{E} |K_{\Delta}(\theta_{k},\alpha_{k})|^{2} \mathbb{E} \|\hat{\nabla}J(\theta_{k})\|^{2} + \left[\operatorname{Cov}\left(|K_{\Delta}(\theta_{k},\alpha_{k})|^{2}, \|\hat{\nabla}J(\theta_{k})\|^{2}\right)\right]$$

26

By Cauchy-Schwarz,

$$\operatorname{Cov}\left(|K_{\Delta}(\theta_{k},\alpha_{k})|^{2},\|\hat{\nabla}J(\theta_{k})\|^{2}\right) \leq \operatorname{Var}\left(|K_{\Delta}(\theta_{k},\alpha_{k})|^{2}\right)\operatorname{Var}\left(\|\hat{\nabla}J(\theta_{k})\|^{2}\right)$$
$$\leq \mathbb{E}\left[|K_{\Delta}(\theta_{k},\alpha_{k})|^{2}\right]\mathbb{E}\left[\|\hat{\nabla}J(\theta_{k})\|^{2}\right] \leq \mathbb{E}|K_{\Delta}(\theta_{k},\alpha_{k})|^{2}\mathbb{E}\|\hat{\nabla}J(\theta_{k})\|^{2}$$

 \mathbf{SO}

$$\mathbb{E} \|K_{\Delta}(\theta_k, \alpha_k) \hat{\nabla} J(\theta_k)\|^2 \le 2\mathbb{E} |K_{\Delta}(\theta_k, \alpha_k)|^2 \mathbb{E} \|\hat{\nabla} J(\theta_k)\|^2$$

Then we bound

$$\begin{split} & \mathbb{E}\left[\|\hat{\nabla}J(\theta_k)\|^2\right] \\ & \leq 3(\mathbb{E}\|L_{\nabla J}|\theta_k| + B\|^2 + \zeta) = 3(\mathbb{E}\left[L_{\nabla J}^2|\theta_k|^2\right] + 2\mathbb{E}\sup_{k\in\mathbb{N}}\left[L_{\nabla J}^2B^2|\theta_k|^2\right] + B^2 + \zeta) \\ & \leq 3(L_{\nabla J}^2(M_\theta + 2B^2M_\theta) + B^2 + \zeta) =: C_0 \end{split}$$

and thus we have

(7.7)
$$\mathbb{E} \|K_{\Delta}(\theta_k, \alpha_k) \hat{\nabla} J(\theta_k)\|^2 \le 6 \mathbb{E} \|K_{\Delta}(\theta_k, \alpha_k)\|^2 \left(L_{\nabla J}^2(M_{\theta} + 2B^2 M_{\theta}) + B^2 + \zeta \right)$$

so we can control this quantity directly by controlling $\mathbb{E} \|K_{\Delta}(\theta_k, \alpha_k)\|^2$, which is done as follows: Notice that

$$\mathbb{P}\left(\frac{1}{\Delta^{2N}}K^2\left(\frac{|\theta_k - \alpha_k|}{\Delta}\right) > x\right) = \mathbb{P}\left(K(\theta_k, \alpha_k) > K(\Delta K^{2^{-1}}(x\Delta^{2N}))\right)$$

By Markov's Inequality we have

$$\mathbb{P}\left(K(\theta_k, \alpha_k) > K(\Delta K^{2^{-1}}(x\Delta^{2N})))\right) \le \frac{\mathbb{E}\|K(\theta_k, \alpha_k)\|}{K(\Delta K^{2^{-1}}(x\Delta^{2N})))} \le \frac{\hat{K}}{K(\Delta K^{2^{-1}}(x\Delta^{2N})))}$$

Then, for all $x \in [\epsilon, \hat{K}_{\epsilon}]$, by choosing Δ as (7.4) we have

$$\mathbb{P}\left(\frac{1}{\Delta^{2N}}K^2\left(\frac{|\theta_k - \alpha_k|}{\Delta}\right) > x\right) \le \frac{2\epsilon}{\sqrt{2\pi}}e^{-x^2}$$

So now observe

$$\begin{split} \mathbb{E} \| K_{\Delta}(\theta_k, \alpha_k) \|^2 &= \mathbb{E} \left[\frac{1}{\Delta^{2N}} K^2 \left(\frac{|\theta_k - \alpha_k|}{\Delta} \right) \right] \\ &= \int_0^\infty \mathbb{P} \left(\frac{1}{\Delta^{2N}} K^2 \left(\frac{|\theta_k - \alpha_k|}{\Delta} \right) > x \right) dx \\ &= \int_0^\epsilon \mathbb{P} \left(\frac{1}{\Delta^{2N}} K^2 \left(\frac{|\theta_k - \alpha_k|}{\Delta} \right) > x \right) dx + \int_\epsilon^{\hat{K}_\epsilon} \mathbb{P} \left(\frac{1}{\Delta^{2N}} K^2 \left(\frac{|\theta_k - \alpha_k|}{\Delta} \right) > x \right) dx \\ &\leq \epsilon + \int_\epsilon^{\hat{K}_\epsilon} \frac{2\epsilon}{\sqrt{2\pi}} e^{-x^2} dx \leq 2\epsilon \end{split}$$

Proof of Lemma 7.7.

Proof.

$$\mathbb{E}\left[\|\pi_{0,\gamma}(\alpha_k)\|^2\right] = \int_0^\infty \mathbb{P}\left(\pi_{0,\gamma}^2(\alpha_k) > x\right) dx$$

We define the level set

$$\Delta_x^{\gamma} := \{ y \in \mathbb{R}^N : \pi_{0,\gamma}^2(y) > x \}$$

so that

$$\mathbb{P}\left(\pi_{0,\gamma}^{2}(\alpha_{k}) > x\right) = \mathbb{P}(\alpha_{k} \in \Delta_{x}^{\gamma})$$

Now split this term as

$$\mathbb{P}(\alpha_{k} \in \Delta_{x}^{\gamma}) = \mathbb{P}(\alpha_{k} \in \Delta_{x}^{\gamma} | \alpha_{k-1} \in \Delta_{x}^{\gamma}) \mathbb{P}(\alpha_{k-1} \in \Delta_{x}^{\gamma}) \\ + \mathbb{P}(\alpha_{k} \in \Delta_{x}^{\gamma} | \alpha_{k-1} \notin \Delta_{x}^{\gamma}) \mathbb{P}(\alpha_{k-1} \notin \Delta_{x}^{\gamma}) \\ \leq \mathbb{P}(\alpha_{k} \in \Delta_{x}^{\gamma} | \alpha_{k-1} \in \Delta_{x}^{\gamma}) + \mathbb{P}(\alpha_{k} \in \Delta_{x}^{\gamma} | \alpha_{k-1} \notin \Delta_{x}^{\gamma})$$

Now, from (2.9), denote

$$\nabla_k := \epsilon \left[K \left(\frac{\theta_k - \alpha_k}{\Delta} \right) \frac{\beta}{2} \hat{\nabla} J(\theta_k) + \nabla \pi_{0,\gamma}(\alpha_k) \right] \pi_{0,\gamma}(\alpha_k), \quad \tilde{w}_k := \sqrt{\epsilon} \pi_{0,\gamma}(\alpha_k) w_k$$

so that $\tilde{w}_k \sim \mathcal{N}(0, \epsilon \pi_{0,\gamma}(\alpha_k)^2)$. Then observe that, given \tilde{w}_k is symmetric with mean zero,

$$\mathbb{P}(\alpha_k \in \Delta_x^{\gamma} | \alpha_{k-1} \in \Delta_x^{\gamma}) = \mathbb{P}(\alpha_{k-1} - \nabla_{k-1} + \tilde{w}_k \in \Delta_x^{\gamma} | \alpha_{k-1} \in \Delta_x^{\gamma})$$
$$\leq \mathbb{P}(\alpha_k \in \Delta_x^{\gamma} | \alpha_{k-1} - \nabla_{k-1} \in \Delta_x^{\gamma})$$

Similarly,

$$\mathbb{P}(\alpha_k \in \Delta_x^{\gamma} | \alpha_{k-1} \notin \Delta_x^{\gamma}) \le \mathbb{P}(\alpha_k \in \Delta_x^{\gamma} | \alpha_{k-1} - \nabla_k \in \Delta_x^{\gamma})$$

so that

$$\mathbb{P}(\alpha_k \in \Delta_x^{\gamma}) \le 2\mathbb{P}(\alpha_k \in \Delta_x^{\gamma} | \alpha_{k-1} - \nabla_{k-1} \in \Delta_x^{\gamma})$$

Now, let $\zeta_z^{\epsilon,\gamma} := (\epsilon \pi_{0,\gamma}^2(z))^{-1}$, define

$$\hat{\Delta}_x^{\gamma}(z) := \{y \in \mathbb{R}^N : y/z \in \Delta_x^{\gamma}\}$$

and notice that

$$\mathbb{P}\left(\alpha_{k} \in \Delta_{x}^{\gamma} | \alpha_{k-1} - \nabla_{k-1} \in \Delta_{x}^{\gamma}\right)$$

$$= \mathbb{P}\left(\zeta_{\alpha_{k-1}}^{\epsilon,\gamma} \alpha_{k} \in \hat{\Delta}_{x}^{\gamma}(\zeta_{\alpha_{k-1}}^{\epsilon,\gamma}) \mid \zeta_{\alpha_{k-1}}^{\epsilon,\gamma} \alpha_{k-1} - \zeta_{\alpha_{k-1}}^{\epsilon,\gamma} \nabla_{k-1} \in \hat{\Delta}_{x}^{\gamma}(\zeta_{\alpha_{k-1}}^{\epsilon,\gamma})\right)$$

$$= \mathbb{P}\left(\zeta_{\alpha_{k-1}}^{\epsilon,\gamma} (\alpha_{k-1} - \nabla_{k-1} + \sqrt{\epsilon}\pi_{0,\gamma}(\alpha_{k-1})w_{k}) \in \hat{\Delta}_{x}^{\gamma}(\zeta_{\alpha_{k-1}}^{\epsilon,\gamma}) \mid \zeta_{\alpha_{k-1}}^{\epsilon,\gamma}(\alpha_{k-1} - \nabla_{k-1}) \in \hat{\Delta}_{x}^{\gamma}(\zeta_{\alpha_{k-1}}^{\epsilon,\gamma})\right)$$

Then, since $\zeta_{\alpha_{k-1}}^{\epsilon,\gamma} \sqrt{\epsilon} \pi_{0,\gamma}(\alpha_{k-1}) w_k \sim \mathcal{N}(0, (\zeta_{\alpha_{k-1}}^{\epsilon,\gamma})^2 \epsilon \pi_{0,\gamma}^2(\alpha_{k-1}))$ we have that,

$$\mathbb{P}\left(\alpha_{k} \in \Delta_{x}^{\gamma} | \alpha_{k-1} - \nabla_{k-1} \in \Delta_{x}^{\gamma}\right) \\
\leq \int_{\hat{\Delta}_{x}^{\gamma}(\zeta_{\alpha_{k-1}}^{\epsilon,\gamma})} \mathcal{N}(\gamma; \hat{c}, (\zeta_{\alpha_{k-1}}^{\epsilon,\gamma})^{2} \epsilon \pi_{0,\gamma}^{2}(\alpha_{k-1})) d\gamma = \int_{\hat{\Delta}_{x}^{\gamma}(\zeta_{\alpha_{k-1}}^{\epsilon,\gamma})} \mathcal{N}(\gamma; \hat{c}, \zeta_{\alpha_{k-1}}^{\epsilon,\gamma}) d\gamma \\
= 28$$

Now, crucially, observe that the volume of $\hat{\Delta}_x^{\gamma}(\zeta_z^{\epsilon,\gamma})$ scales, w.r.t z, at the same rate as the variance of $\mathcal{N}(\cdot, \hat{c}, \zeta_z^{\epsilon,\gamma})$. Thus we have

$$\int_{\hat{\Delta}_x^{\gamma}(\zeta_{z_1}^{\epsilon,\gamma})} \mathcal{N}(\gamma;\hat{c},\zeta_{z_1}^{\epsilon,\gamma}) d\gamma = \int_{\hat{\Delta}_x^{\gamma}(\zeta_{z_2}^{\epsilon,\gamma})} \mathcal{N}(\gamma;\hat{c},\zeta_{z_2}^{\epsilon,\gamma}) d\gamma \quad \forall z_1, z_2 > 0$$

In particular, take z such that $\zeta_z^{\epsilon,\gamma} = (\epsilon \pi_{0,\gamma}^2(z))^{-1} = \epsilon^2$. Note that this necessitates $\bar{\pi}_{0,\gamma} \geq (\frac{1}{\epsilon})^{3/2}$, which is given from the condition $\gamma \leq \bar{\pi}_0 \epsilon^{3/2}$. Then we have

$$\mathbb{P}\left(\alpha_{k} \in \Delta_{x}^{\gamma} | \alpha_{k-1} - \nabla_{k-1} \in \Delta_{x}^{\gamma}\right) \leq \int_{\hat{\Delta}_{x}^{\gamma}(\epsilon^{2})} \mathcal{N}(\gamma; \hat{c}, \epsilon^{2}) d\gamma$$
$$\leq \int_{\hat{\Delta}_{x}^{\gamma}(\epsilon^{2})} \frac{1}{\epsilon \sqrt{2\pi}} d\gamma = \int_{\hat{\Delta}_{x}^{\gamma}(1)} \frac{\epsilon}{\sqrt{2\pi}} d\gamma$$

Then

$$\mathbb{E}\left[\|\pi_{0,\gamma}(\alpha_k)\|^2\right] = \int_0^\infty \mathbb{P}(\pi_{0,\gamma}^2(\alpha_k) > x)dx$$

$$\leq 2\int_0^\infty \int_{\Delta_x^\gamma} \frac{\epsilon}{\sqrt{2\pi}} d\gamma dx = \sqrt{\frac{2}{\pi}}\epsilon \int_0^\infty \int_{\mathbb{R}^N} \mathbb{1}\{\pi_{0,\gamma}(\gamma) > \sqrt{x}\}d\gamma dx$$

but observe that, for $\gamma \leq 1$,

$$\int_0^\infty \int_{\mathbb{R}^N} \mathbb{1}\{\pi_{0,\gamma}(\gamma) > \sqrt{x}\} d\gamma \, dx \le \int_0^\infty \int_{\mathbb{R}^N} \mathbb{1}\{\pi_0(\gamma) > \sqrt{x}\} d\gamma \, dx$$

So now

$$\mathbb{E}\left[\|\pi_{0,\gamma}(\alpha_{k})\|^{2}\right] \leq \sqrt{\frac{2}{\pi}} \epsilon \int_{0}^{\infty} \int_{\mathbb{R}^{N}} \mathbb{1}\{\pi_{0}(\gamma) > \sqrt{x}\} d\gamma dx$$
$$= \sqrt{\frac{2}{\pi}} \epsilon \left[\int_{0}^{1} \int_{\mathbb{R}^{N}} \mathbb{1}\{\pi_{0}(\gamma) > \sqrt{x}\} d\gamma dx$$
$$+ \int_{1}^{\overline{\pi}_{0,\gamma}^{2}} \int_{\mathbb{R}^{N}} \mathbb{1}\{\pi_{0}(\gamma) > \sqrt{x}\} d\gamma dx\right]$$
$$\leq \sqrt{\frac{2}{\pi}} \epsilon \left[\int_{0}^{1} \int_{\mathbb{R}^{N}} \mathbb{1}\{\pi_{0}(\gamma) > x\} d\gamma dx + \overline{\pi}_{0,\gamma}^{2} V(\{\pi_{0} > 1\})\right]$$
$$\leq \epsilon \left[1 + V(\{\pi_{0} > 1\})\right] =: \epsilon V_{1}$$

where $V({\pi_0 > 1})$ is shorthand for $\int_{\mathbb{R}^N} \mathbb{1}{\pi_0(\gamma) > 1}d\gamma$.

Define, analagously,

$$\Gamma_x^{\gamma} = \{ y \in \mathbb{R}^N : \| \nabla \pi_{0,\gamma}(y) \|^2 > x \}$$

By the same procedure as above, we can obtain

$$\mathbb{P}(\alpha_k \in \Gamma_x^{\gamma}) \le 2 \int_{\Gamma_x^{\gamma}} \frac{\epsilon}{\sqrt{2\pi}} d\gamma$$

and so

$$\mathbb{E} \|\nabla \pi_{0,\gamma}(\alpha_k)\|^2 = \int_0^\infty \mathbb{P}(\|\nabla \pi_{0,\gamma}(\alpha_k)\|^2 > x) dx \le 2 \int_0^\infty \int_{\Gamma_x^\gamma} \sqrt{\frac{2}{\pi}} \epsilon d\gamma dx$$
$$\le \epsilon \left[1 + \bar{\pi}'_0 V(\{\|\nabla \pi_0\| > 1\})\right] =: \epsilon V_2$$
29

where $\bar{\pi}'_0 = \sup_{x \in \mathbb{R}^N} \nabla \pi_0(x)$, $V(\{\|\nabla \pi_0\| > 1\}) = \int_{\mathbb{R}^N} \mathbb{1}\{\|\nabla \pi_0(\gamma)\| > 1\} d\gamma$ Taking π_0 such that $\bar{\pi}_0 = 1$ gives

$$\mathbb{E}\left[\|\pi_{0,\gamma}(\alpha_k)\|^2\right] \le \epsilon, \quad \mathbb{E}\|\nabla\pi_{0,\gamma}(\alpha_k)\|^2 \le \epsilon \qquad \Box$$

Proof of Lemma 7.8.

Proof. First take

$$\mathbb{E} \|\pi_{0,\gamma}(\bar{\alpha}(s))R\|^{2} = \mathbb{E} \left[|\pi_{0,\gamma}(\bar{\alpha}(s))|^{2} \|R\|^{2} \right] \\ = \mathbb{E} \left[|\pi_{0,\gamma}(\bar{\alpha}(s))|^{2} \right] \mathbb{E} \left[\|R\|^{2} \right] + \operatorname{Cov} \left(|\pi_{0,\gamma}(\bar{\alpha}(s))|^{2}, \|R\|^{2} \right)$$

By Cauchy-Schwarz,

$$\operatorname{Cov}\left(|\pi_{0,\gamma}(\bar{\alpha}(s))|^{2}, \|R\|^{2}\right) \leq \operatorname{Var}\left(|\pi_{0,\gamma}(\bar{\alpha}(s))|^{2}\right) \operatorname{Var}\left(\|R\|^{2}\right)$$
$$\leq \mathbb{E}\left[|\pi_{0,\gamma}(\bar{\alpha}(s))|^{2}\right] \mathbb{E}\left[\|R\|^{2}\right]$$

So we have

(7.8)
$$\mathbb{E}\|\pi_{0,\gamma}(\bar{\alpha}(s))R\|^2 \le 2\mathbb{E}\left[|\pi_{0,\gamma}(\bar{\alpha}(s))|^2\right]\mathbb{E}\left[\|R\|^2\right] \qquad \Box$$

Proof of Lemma 7.9.

Proof. This follows from Assumption 3.4

Proof of Lemma 7.10.

Proof. Lemma 7.10 equivalent to Lemma 2 of [22], and follows from Assumptions A1, A2, and A3. $\hfill \Box$

Proof of Lemma 7.11.

Proof. First note that for all $x \in \mathbb{R}^N$ we have $\nabla \pi_{0,\gamma}(x) = \frac{\nabla \pi_0(x)}{\gamma}$, so

$$\arg\max_{y\in\mathbb{R}^N} \|\langle y, \, \nabla\pi_{0,\gamma}(y)\rangle\| = \arg\max_{y\in\mathbb{R}^N} \|\langle\gamma y, \, \frac{\nabla\pi_0(y)}{\gamma}\rangle\| = \arg\max_{y\in\mathbb{R}^N} \|\langle y, \, \nabla\pi_0(y)\rangle\|$$

and

$$\arg \max_{y \in \mathbb{R}^N} \| \langle y, \pi_{0,\gamma}(y) \nabla \pi_{0,\gamma}(y) \rangle \| = \arg \max_{y \in \mathbb{R}^N} \| \langle \gamma y, \pi_{0,\gamma}(y) \frac{\nabla \pi_0(y)}{\gamma} \rangle \|$$
$$= \arg \max_{y \in \mathbb{R}^N} \| \langle y, \pi_0(y) \nabla \pi_0(y) \rangle \|$$

Thus (7.6) is equivalent to:

(7.9)
$$\exists I > 0 : ||\langle x, \nabla \pi_0(x) \rangle|| \le I \ \forall x \in \mathbb{R}^N$$

(7.10)
$$\exists I' > 0 : ||\langle x, \pi_0(x) \nabla \pi_0(x) \rangle|| \le I' \ \forall x \in \mathbb{R}^N$$

Now we prove by reductio ad absurdum: Suppose (7.9), (7.10) do not hold. Then we have:

(7.11)
$$\begin{aligned} \forall y > 0 \ \exists x \in \mathbb{R}^N : \ \|\langle x, \nabla \pi_0(x) \rangle\| > y \\ \forall y > 0 \ \exists x \in \mathbb{R}^N : \ \|\langle x, \pi_0(x) \nabla \pi_0(x) \rangle\| > y \\ 30 \end{aligned}$$

This manuscript is for review purposes only.

Recall M as defined in Assumption 3.4. Denote

$$D_1 = \{ x \in \mathbb{R}^N : \|x\| \le M \lor 1 \}, \ D_2 = D_1^C = \mathbb{R}^N \backslash D_1$$

Recall that we assume Lipshitz-continuity of $\pi_0(\cdot)$ in (A3.5), with Lipschitz constant L_{π_0} . Thus, $\|\nabla \pi_0(x)\|$ is bounded by L_{π_0} for all $x \in \mathbb{R}^N$. In particular notice that since $\|\nabla \pi_0(x)\|$, $\|\pi_0(x)\|$, and $\|x\|$ are bounded for $x \in D_1$, there exists some M^*, M^{**} such that

$$\|\langle x, \nabla \pi_0(x) \rangle\| \le M^* \ \forall x \in D_1$$

$$|\langle x, \pi_0(x) \nabla \pi_0(x) \rangle|| \le M^{**} \ \forall x \in D_1$$

Then (7.11) requires both:

(7.12)
$$\begin{aligned} \forall y > M^* \; \exists x \in D_2 : \; \|\langle x, \nabla \pi_0(x) \rangle\| > y \\ \forall y > M^{**} \; \exists x \in D_2 : \; \|\langle x, \nabla \pi_0(x) \nabla \pi_0(x) \rangle\| > y \end{aligned}$$

But by assumption 7.9 $\nabla \pi_0(x)$ decays as $\mathcal{O}(||x||^{-1})$ for $x \in D_2$, and in Lemma 7.9 $\pi_0(x)$ decays as $\mathcal{O}(||x||^{-2})$. Specifically, there exists \tilde{C} such that

$$\|\nabla \pi_0(x)\| \le \frac{\tilde{C}}{\|x\|^{-1}}, \quad \|\pi_0(x)\| \le \frac{2}{(\beta m^*)^2 \|x\|^2}, \, \forall x \in D_2$$

Thus we have that

$$\|\langle x, \nabla \pi_0(x) \rangle\| \le \tilde{C} \ \forall x \in D_2$$
$$\|\langle x, \pi_0(x) \nabla \pi_0(x) \rangle\| \le \frac{2\tilde{C}}{(\beta m^*)^2} \ \forall x \in D_2$$

which contradicts (7.12) and thus (7.11) is refuted. So (7.6) holds.

7.4. Proofs of Section 5 Results.

Proof of Lemma 5.1.

Proof. We begin by bounding this Wasserstein distance by the mean-square error between the processes (5.2) and (5.3). Recall that $Y(k\epsilon)$ has probability law π_k .

$$\mathcal{W}_2(\pi_k, \gamma_{k\epsilon}) = \inf_{\gamma \in \Gamma(\pi_k, \gamma_{k\epsilon})} \left(\mathbb{E}_{(x,y) \sim \gamma} \|x - y\|^2 \right)^{1/2} \\ \leq \sqrt{\mathbb{E}_{x \sim \pi_k, y \sim \gamma_{k\epsilon}} \|x - y\|^2}$$

Then we take $t = k\epsilon$, and bound

(7.13)
$$\mathbb{E} \|Y(t) - X(t)\|^{2} \leq 3\mathbb{E} \|\int_{0}^{t} g_{s}(\theta_{\bar{s}}, Y(s)) - \hat{g}_{s}(\theta_{\bar{s}}, X(s)) ds\|^{2} + 3\mathbb{E} \|\int_{0}^{t} \mathbb{E} [\pi_{0,\gamma}(\bar{\alpha}(\bar{s}))|\bar{\alpha}(s) = Y(s)] - \pi_{0,\gamma}(X(s)) dW(s)\|^{2}$$

First we bound, using Jensen's inequality:

$$\mathbb{E} \| \int_0^t g_s(\theta_{\bar{s}}, Y(s)) - \hat{g}_s(\theta_{\bar{s}}, X(s)) ds \|^2$$
31

$$\begin{split} &= \mathbb{E} \| \int_{0}^{t} \mathbb{E} \left[\left(K_{\Delta}(\theta_{\bar{s}}, \bar{\alpha}(\bar{s})) \frac{\beta}{2} \hat{\nabla} J(\theta_{\bar{s}}) + \nabla \pi_{0,\gamma}(\bar{\alpha}(\bar{s})) \right) \pi_{0,\gamma}(\bar{\alpha}(\bar{s})) | \bar{\alpha}(s) = Y(s) \right] \\ &\quad - \mathbb{E} \left[\left(K_{\Delta}(\theta_{\bar{s}}, \bar{\alpha}(\bar{s})) \frac{\beta}{2} \hat{\nabla} J(\theta_{\bar{s}}) + \nabla \pi_{0,\gamma}(\bar{\alpha}(\bar{s})) \right) \pi_{0,\gamma}(\bar{\alpha}(\bar{s})) | \bar{\alpha}(s) = X(s) \right] ds \|^{2} \\ &\leq \mathbb{E} \| \int_{0}^{t} \sqrt{\mathbb{E} \left[\| K_{\Delta}(\theta_{\bar{s}}, \bar{\alpha}(\bar{s})) \frac{\beta}{2} \hat{\nabla} J(\theta_{\bar{s}}) \pi_{0,\gamma}(\bar{\alpha}(\bar{s})) \|^{2} | \bar{\alpha}(s) = Y(s) \right]} \\ &\quad - \sqrt{\mathbb{E} \left[\| K_{\Delta}(\theta_{\bar{s}}, \bar{\alpha}(\bar{s})) \frac{\beta}{2} \hat{\nabla} J(\theta_{\bar{s}}) \pi_{0,\gamma}(\bar{\alpha}(\bar{s})) \|^{2} | \bar{\alpha}(s) = X(s) \right]} \\ &\quad + \sqrt{\mathbb{E} \left[\| \pi_{0,\gamma}(\bar{\alpha}(\bar{s})) \nabla \pi_{0,\gamma}(\bar{\alpha}(\bar{s})) \|^{2} | \bar{\alpha}(s) = Y(s) \right]} \\ &\quad - \sqrt{\mathbb{E} \left[\| \pi_{0,\gamma}(\bar{\alpha}(\bar{s})) \nabla \pi_{0,\gamma}(\bar{\alpha}(\bar{s})) \|^{2} | \bar{\alpha}(s) = X(s) \right]} ds \|^{2} \\ &\leq \mathbb{E} \| \int_{0}^{t} 4 \epsilon \sqrt{2C_{0}} + 2\sqrt{2} \epsilon ds \|^{2} \leq (k\epsilon)^{2} \epsilon^{2} (96C_{0} + 24) \end{split}$$

We bound the second term, using the Itô Isometry:

$$\mathbb{E} \| \int_0^t \mathbb{E} \left[\pi_{0,\gamma}(\bar{\alpha}(\bar{s})) | \bar{\alpha}(s) = Y(s) \right] - \pi_{0,\gamma}(X(s)) dW(s) \|^2$$
$$= \mathbb{E} \| \int_0^t \mathbb{E} \left[\pi_{0,\gamma}(\bar{\alpha}(\bar{s})) - \pi_{0,\gamma}(X(s)) | \bar{\alpha}(s) = Y(s) \right] dW(s) \|^2$$
$$= \mathbb{E} \left[\int_0^t \mathbb{E} \left[\| \pi_{0,\gamma}(\bar{\alpha}(\bar{s})) - \pi_{0,\gamma}(X(s)) \|^2 | \bar{\alpha}(s) = Y(s) \right] ds \right]$$
$$\leq 6 \mathbb{E} \left[\int_0^t \epsilon ds \right] = 6 k \epsilon^2$$

Thus, we have

$$\mathcal{W}_2(\pi_k, \gamma_{k\epsilon}) \le \sqrt{72(k\epsilon)^2 \epsilon^2 (4C_0 + 1) + 18(k\epsilon)\epsilon} \\ \le 6 (k\epsilon)\epsilon \sqrt{12C_0 + 3} + 3\sqrt{2 (k\epsilon)\epsilon} \qquad \Box$$

Proof of Lemma 5.2.

Proof. Let $\mathbb{P}_X^t := \text{Law}(X(s) : 0 \le s \le t)$ for X(s) in 5.3 and $\mathbb{P}_A^t := \text{Law}(\alpha(s) : 0 \le s \le t)$ for $\alpha(t)$ in 2.12. The Radon-Nikodym derivative of \mathbb{P}_A^t with respect to \mathbb{P}_X^t is given by the Girsanov formula:

$$\begin{split} \frac{d\mathbb{P}_{A}^{t}}{d\mathbb{P}_{X}^{t}}(X) &= \exp\{\int_{0}^{t} \left(G(X(s)) - g_{s}(\theta_{\bar{s}}, X(s))\right)^{*} \pi_{0,\gamma}(X(s))^{-2} dW(s) \\ &- \frac{1}{2} \int_{0}^{t} \|G(X(s)) - g_{s}(\theta_{\bar{s}}, X(s)))\|^{2} \pi_{0,\gamma}(X(s))^{-2} ds\} \end{split}$$

where $G(X(s)) = \frac{\beta}{2} \pi_{0,\gamma}^2(X(s)) \nabla J(X(s)) + \pi_{0,\gamma}(X(s)) \nabla \pi_{0,\gamma}(X(s))$ Then by the martingale property of the Itô integral,

$$D(\mathbb{P}_X^t \| \mathbb{P}_A^t) = -\int d\mathbb{P}_X^t \log \frac{d\mathbb{P}_A^t}{d\mathbb{P}_X^t}$$
32

$$\begin{split} &= \frac{1}{2} \int_0^t \mathbb{E} \bigg[\|G(X(s)) - g_s(\theta_{\bar{s}}, X(s)))\|^2 \pi_{0,\gamma}(X(s))^{-2} \bigg] ds \\ &= \frac{1}{2} \int_0^t \mathbb{E} \bigg[\|\frac{\beta}{2} \pi_{0,\gamma}^2(X(s)) \nabla J(X(s)) + \pi_{0,\gamma}(X(s)) \nabla \pi_{0,\gamma}(X(s)) \\ &- g_s(\theta_{\bar{s}}, X(s))\|^2 \pi_{0,\gamma}(X(s))^{-2} \bigg] ds \\ &= \frac{1}{2} \int_0^t \mathbb{E} \bigg[\|\frac{\beta}{2} \pi_{0,\gamma}^2(\bar{\alpha}(s)) \nabla J(\bar{\alpha}(s)) + \pi_{0,\gamma}(\bar{\alpha}(s)) \nabla \pi_{0,\gamma}(\bar{\alpha}(s)) \\ &- g_s(\theta_{\bar{s}}, \bar{\alpha}(s))\|^2 \pi_{0,\gamma}(\bar{\alpha}(s))^{-2} \bigg] ds \end{split}$$

where the last line follows from the fact that $\text{Law}(\bar{\alpha}(s)) = \text{Law}(X(s)) \ \forall s \leq t$.

Now let $t = k\epsilon$ for some $k \in \mathbb{N}$. Then, expanding $g_s(\theta_{\bar{s}}, Y(s))$ and using Jensen's inequality, we get

$$\begin{aligned} (7.14) \\ D(\mathbb{P}_X^t \| \mathbb{P}_A^t) &\leq \frac{1}{2} \int_0^t \mathbb{E} \left[\| \frac{\beta}{2} \pi_{0,\gamma}^2(\bar{\alpha}(s)) \nabla J(\bar{\alpha}(s)) + \pi_{0,\gamma}(\bar{\alpha}(s)) \nabla \pi_{0,\gamma}(\bar{\alpha}(s)) \\ &- \left(K_{\Delta}(\theta_{\bar{s}}, \bar{\alpha}(s)) \frac{\beta}{2} \hat{\nabla} J(\theta_{\bar{s}}) + \nabla \pi_{0,\gamma}(\bar{\alpha}(s)) \right) \pi_{0,\gamma}(\bar{\alpha}(s)) \|^2 \pi_{0,\gamma}(\bar{\alpha}(s))^{-2} \right] ds \\ &= \frac{\beta}{4} \int_0^t \mathbb{E} \| \pi_{0,\gamma}(\bar{\alpha}(s)) \nabla J(\bar{\alpha}(s)) - K_{\Delta}(\theta_{\bar{s}}, \bar{\alpha}(s)) \hat{\nabla} J(\theta_{\bar{s}}) \|^2 ds \\ &= \frac{\beta}{4} \sum_{j=0}^{k-1} \int_{j\epsilon}^{(j+1)\epsilon} \mathbb{E} \| \pi_{0,\gamma}(\bar{\alpha}(s)) \nabla J(\bar{\alpha}(s)) - K_{\Delta}(\theta_{\bar{s}}, \bar{\alpha}(s)) \hat{\nabla} J(\theta_{\bar{s}}) \|^2 ds \\ &\leq \frac{\beta}{2} \sum_{j=0}^{k-1} \int_{j\epsilon}^{(j+1)\epsilon} \mathbb{E} \| \pi_{0,\gamma}(\bar{\alpha}(s)) \nabla J(\bar{\alpha}(s)) - \pi_{0,\gamma}(\bar{\alpha}(s)) \nabla J(\bar{\alpha}(\lfloor s/\epsilon \rfloor \epsilon)) \|^2 ds \\ &+ \frac{\beta}{2} \sum_{j=0}^{k-1} \int_{j\epsilon}^{(j+1)\epsilon} \mathbb{E} \| \pi_{0,\gamma}(\bar{\alpha}(s)) |^2 \mathbb{E} \| \nabla J(\bar{\alpha}(\lfloor s/\epsilon \rfloor \epsilon)) - K_{\Delta}(\theta_{\bar{s}}, \bar{\alpha}(s)) \hat{\nabla} J(\theta_{\bar{s}}) \|^2 ds \\ &\leq \beta \sum_{j=0}^{k-1} \int_{j\epsilon}^{(j+1)\epsilon} \mathbb{E} \| \pi_{0,\gamma}(\bar{\alpha}(s)) |^2 \mathbb{E} \| \nabla J(\bar{\alpha}(\lfloor s/\epsilon \rfloor \epsilon)) - \nabla J(\alpha(0)) \|^2 \\ &+ 3\beta \sum_{j=0}^{k-1} \int_{j\epsilon}^{(j+1)\epsilon} \mathbb{E} \| \pi_{0,\gamma}(\bar{\alpha}(s)) |^2 \mathbb{E} \| \nabla J(\bar{\alpha}(\lfloor s/\epsilon \rfloor \epsilon)) - \nabla J(\alpha(0)) \|^2 \\ &+ 3\beta \sum_{j=0}^{k-1} \int_{j\epsilon}^{(j+1)\epsilon} \mathbb{E} \| \pi_{0,\gamma}(\bar{\alpha}(s)) |^2 \mathbb{E} \| \nabla J(\bar{\alpha}(0)) \|^2 ds \end{aligned}$$

where we use Lemma 7.8 in the last inequality.

First Term: By $L_{\nabla J}$ -smoothness (Assumption 3.1), we begin to control the first

term as:

$$\mathbb{E} \|\nabla J(\bar{\alpha}(s)) - \nabla J(\bar{\alpha}(\lfloor s/\epsilon \rfloor \epsilon))\|^2 \le L^2_{\nabla J} \mathbb{E} \|(\bar{\alpha}(s) - \bar{\alpha}(\lfloor s/\epsilon \rfloor \epsilon))\|^2$$

Then for $s \in [j\epsilon, (j+1)\epsilon]$:

$$\begin{aligned} &(\bar{\alpha}(s) - \bar{\alpha}(\lfloor s/\epsilon \rfloor \epsilon)) \\ &= (s - j\epsilon) \frac{\beta}{2} K_{\Delta}(\theta_j, \bar{\alpha}(j\epsilon)) \hat{\nabla} J(\theta_j) \pi_{0,\gamma}(\bar{\alpha}(j\epsilon)) \\ &+ (s - j\epsilon) \pi_{0,\gamma}(\bar{\alpha}(j\epsilon)) \nabla \pi_{0,\gamma}(\bar{\alpha}(j\epsilon)) - \pi_{0,\gamma}(\bar{\alpha}(j\epsilon)) (W(s) - W(j\epsilon)) \end{aligned}$$

Then bound $\mathbb{E} \| (\bar{\alpha}(s) - \bar{\alpha}(\lfloor s/\epsilon \rfloor \epsilon)) \|^2$ as:

$$\begin{aligned}
\mathbb{E} \| (\bar{\alpha}(s) - \bar{\alpha}(\lfloor s/\epsilon \rfloor \epsilon)) \|^2 \\
(7.15) &\leq 6\epsilon^2 \mathbb{E} \| K_{\Delta}(\theta_j, \bar{\alpha}(j\epsilon)) \frac{\beta}{2} \hat{\nabla} J(\theta_j) \|^2 \mathbb{E} \| \pi_{0,\gamma}(\bar{\alpha}(j\epsilon)) \|^2 \\
&\quad + 3\epsilon^2 \mathbb{E} \| \pi_{0,\gamma}(\bar{\alpha}(j\epsilon)) \nabla \pi_{0,\gamma}(\bar{\alpha}(j\epsilon)) \|^2 + 3\epsilon \mathbb{E} \| \pi_{0,\gamma}(\bar{\alpha}(j\epsilon)) (W(s) - W(j\epsilon) \|^2
\end{aligned}$$

But observe that

$$\mathbb{E}\|(\bar{\alpha}(s) - \bar{\alpha}(\lfloor s/\epsilon \rfloor \epsilon))\|^{2} = \mathbb{E}\left[(s - j\epsilon)^{2}\|K_{\Delta}(\theta_{j}, \bar{\alpha}(j\epsilon))\hat{\nabla}J(\theta_{j})\pi_{0,\gamma}(\bar{\alpha}(j\epsilon))\|^{2}\right] \\ + \mathbb{E}\left[(s - j\epsilon)^{2}\|\pi_{0,\gamma}(\bar{\alpha}(j\epsilon))\nabla\pi_{0,\gamma}(\bar{\alpha}(j\epsilon))\|^{2}\right] \\ + 2\mathbb{E}\left[(s - j\epsilon)^{2}\langle K_{\Delta}(\theta_{j}, \bar{\alpha}(j\epsilon))\hat{\nabla}J(\theta_{j}), \pi_{0,\gamma}^{2}(\bar{\alpha}(j\epsilon))\nabla\pi_{0,\gamma}(\bar{\alpha}(j\epsilon))\rangle\right] \\ - \mathbb{E}\left[(s - j\epsilon)\langle K_{\Delta}(\theta_{j}), \bar{\alpha}(j\epsilon))\hat{\nabla}J(\theta_{j}), \pi_{0,\gamma}^{2}(W(s) - W(j\epsilon))\rangle\right] \\ - \mathbb{E}\left[(s - j\epsilon)\langle \nabla\pi_{0,\gamma}(\bar{\alpha}(j\epsilon)), \pi_{0,\gamma}^{2}(W(s) - W(j\epsilon))\rangle\right] \\ + \mathbb{E}\left[\pi_{0,\gamma}^{2}\|\bar{\alpha}(j\epsilon))(W(s) - W(j\epsilon)\|^{2}\right]$$

then combining (7.16) with (7.15), and using the Martingale property of Brownian motion, Jensen's Inequality, and Cauchy-Schwarz, gives:

$$\begin{split} & \mathbb{E}\left[\pi_{0,\gamma}^{2}(\bar{\alpha}(j\epsilon))\|W(s) - W(j\epsilon)\|^{2}\right] \\ &\leq 2\epsilon^{2}\mathbb{E}\|K_{\Delta}(\theta_{j},\bar{\alpha}(j\epsilon))\hat{\nabla}J(\theta_{j})\pi_{0,\gamma}(\bar{\alpha}(j\epsilon))\|^{2} \\ &\quad + 2\epsilon^{2}\mathbb{E}\|\pi_{0,\gamma}(\bar{\alpha}(j\epsilon))\nabla\pi_{0,\gamma}(\bar{\alpha}(j\epsilon))\|^{2} \\ &\quad + \epsilon^{2}\sqrt{\mathbb{E}}\left[\|\langle K_{\Delta}(\theta_{j},\bar{\alpha}(j\epsilon))\hat{\nabla}J(\theta_{j})\pi_{0,\gamma}^{2}(\bar{\alpha}(j\epsilon)), \pi_{0,\gamma}^{2}(\bar{\alpha}(j\epsilon))\nabla\pi_{0,\gamma}(\bar{\alpha}(j\epsilon))\rangle\|^{2}\right] \\ &\leq 2\epsilon^{2}\mathbb{E}\|K_{\Delta}(\theta_{j},\bar{\alpha}(j\epsilon))\hat{\nabla}J(\theta_{j})\pi_{0,\gamma}(\bar{\alpha}(j\epsilon))\|^{2} \\ &\quad + 2\epsilon^{2}\mathbb{E}\|\pi_{0,\gamma}(\bar{\alpha}(j\epsilon))\nabla\pi_{0,\gamma}(\bar{\alpha}(j\epsilon))\|^{2} \\ &\quad + \epsilon^{2}\sqrt{\mathbb{E}}\|K_{\Delta}(\theta_{j},\bar{\alpha}(j\epsilon))\hat{\nabla}J(\theta_{j})\pi_{0,\gamma}(\bar{\alpha}(j\epsilon))\|^{2} \\ &\quad + \epsilon^{2}\sqrt{\mathbb{E}}\|\pi_{0,\gamma}(\bar{\alpha}(j\epsilon))\nabla\pi_{0,\gamma}(\bar{\alpha}(j\epsilon))\|^{2} \end{split}$$

Thus we have:

$$\mathbb{E}\|(\bar{\alpha}(s) - \bar{\alpha}(\lfloor s/\epsilon \rfloor \epsilon))\|^2$$
34

$$\leq 18\epsilon^{2}\mathbb{E}\|K_{\Delta}(\theta_{j},\bar{\alpha}(j\epsilon))\frac{\beta}{2}\hat{\nabla}J(\theta_{j})\|^{2}\mathbb{E}\|\pi_{0,\gamma}(\bar{\alpha}(j\epsilon))\|^{2}$$

$$+9\epsilon^{2}\mathbb{E}\|\pi_{0,\gamma}(\bar{\alpha}(j\epsilon))\nabla\pi_{0,\gamma}(\bar{\alpha}(j\epsilon))\|^{2}$$

$$+3\epsilon^{2}\sqrt{\mathbb{E}}\|K_{\Delta}(\theta_{j},\bar{\alpha}(j\epsilon))\frac{\beta}{2}\hat{\nabla}J(\theta_{j})\|^{2}\mathbb{E}}\|\pi_{0,\gamma}(\bar{\alpha}(j\epsilon))\|^{2}$$

$$+3\epsilon^{2}\sqrt{\mathbb{E}}\|\pi_{0,\gamma}(\bar{\alpha}(j\epsilon))\nabla\pi_{0,\gamma}(\bar{\alpha}(j\epsilon))\|^{2}$$

$$\leq \epsilon^{4}\left(72C_{0}+18\right)+\epsilon^{3}\left(6\sqrt{C_{0}}+\sqrt{2}\right)$$

where recall $C_0 := 3L_{\nabla J}^2(M_\theta + 2B^2M_\theta) + B^2 + \zeta$. Consequently,

(7.17)
$$\beta \sum_{j=0}^{k-1} \int_{j\epsilon}^{(j+1)\epsilon} \mathbb{E} |\pi_{0,\gamma}(\bar{\alpha}(s))|^2 \mathbb{E} ||\nabla J(\bar{\alpha}(s)) - \nabla J(\bar{\alpha}(\lfloor s/\epsilon \rfloor \epsilon))||^2 ds$$
$$\leq \beta \sum_{j=0}^{k-1} \int_{j\epsilon}^{(j+1)\epsilon} L^2_{\nabla J} \mathbb{E} |\pi_{0,\gamma}(\bar{\alpha}(s))|^2 \mathbb{E} ||(\bar{\alpha}(s) - \bar{\alpha}(\lfloor s/\epsilon \rfloor \epsilon))||^2 ds$$
$$\leq \beta L^2_{\nabla J} k\epsilon \left(\epsilon^4 \left(72C_0 + 18\right) + \epsilon^3 \left(6\sqrt{C_0} + \sqrt{2}\right)\right)$$

 ${\bf Second \ Term:}$ We now bound the second term as

$$\begin{aligned} 3\beta \sum_{j=0}^{k-1} \int_{j\epsilon}^{(j+1)\epsilon} \mathbb{E} |\pi_{0,\gamma}(\bar{\alpha}(s))|^2 \mathbb{E} ||\nabla J(\bar{\alpha}(\lfloor s/\epsilon \rfloor \epsilon)) - \nabla J(\bar{\alpha}(0))||^2 ds \\ &= 3\beta \sum_{j=0}^{k-1} \epsilon \mathbb{E} |\pi_{0,\gamma}(\bar{\alpha}(s))|^2 \mathbb{E} ||\nabla J(\bar{\alpha}(j\epsilon)) - \nabla J(\bar{\alpha}(0))||^2 ds \\ &\leq 3\beta \sum_{j=0}^{k-1} \epsilon \mathbb{E} |\pi_{0,\gamma}(\bar{\alpha}(s))|^2 \mathbb{E} ||\sum_{i=0}^{j-1} |\nabla J(\bar{\alpha}((i+1)\epsilon)) - \nabla J(\bar{\alpha}(i\epsilon))|||^2 \\ &\leq 3\beta \sum_{j=0}^{k-1} \epsilon^2 L_{\nabla J}^2 \mathbb{E} ||\sum_{i=0}^{j-1} |\bar{\alpha}((i+1)\epsilon) - \alpha(i\epsilon)||^2 \\ (7.18) &= 3\beta \sum_{j=0}^{k-1} \epsilon^2 L_{\nabla J}^2 \mathbb{E} \left[\sum_{i=0}^{j-1} \sum_{l=0}^{j-1} |\bar{\alpha}((i+1)\epsilon) - \bar{\alpha}(i\epsilon)||\bar{\alpha}((l+1)\epsilon) - \bar{\alpha}(l\epsilon)|| \right] \\ &\leq 3\beta \sum_{j=0}^{k-1} \epsilon^2 L_{\nabla J}^2 \left[\sum_{i=0}^{j-1} \sum_{l=0}^{j-1} \mathbb{E} \left[|\bar{\alpha}((i+1)\epsilon) - \bar{\alpha}(i\epsilon)||\bar{\alpha}((l+1)\epsilon) - \bar{\alpha}(l\epsilon)|| \right] \\ &\leq 3\beta \sum_{j=0}^{k-1} \epsilon^2 L_{\nabla J}^2 \left[\sum_{i=0}^{j-1} \sum_{l=0}^{j-1} (\epsilon^4 (72C_0 + 18) + \epsilon^3 \left(6\sqrt{C_0} + \sqrt{2}\right)) \right] \\ &= 3\beta k\epsilon^2 L_{\nabla J}^2 k^2 \left(\epsilon^4 (72C_0 + 18) + \epsilon^3 \left(6\sqrt{C_0} + \sqrt{2}\right) \right) \\ &\leq 3\beta k\epsilon^2 L_{\nabla J}^2 \left(\epsilon^4 (72C_0 + 18) + \epsilon^3 \left(6\sqrt{C_0} + \sqrt{2}\right) \right) \\ &= 3\beta (k\epsilon)^3 L_{\nabla J}^2 \left(\epsilon^4 (72C_0 + 18) + \epsilon^3 \left(6\sqrt{C_0} + \sqrt{2}\right) \right) \end{aligned}$$

Third Term: We bound $\mathbb{E} \| \pi_{0,\gamma}(\bar{\alpha}(s)) \nabla(J(\alpha(0))) \|^2$ by controlling the sampling distribution $\pi_{0,\gamma}(\cdot)$. By Lemma 7.8 and Assumption 3.1 we have

(7.19)
$$\mathbb{E}\|\pi_{0,\gamma}(\bar{\alpha}(s))\nabla(J(\alpha(0)))\|^{2} \leq 2\mathbb{E}\left[\|\pi_{0,\gamma}(\bar{\alpha}(s))\|^{2}\right]\mathbb{E}\left[\|\nabla(J(\alpha(0)))\|^{2}\right] \\ \leq 2L_{J}^{2}\mathbb{E}\left[\|\pi_{0,\gamma}(\bar{\alpha}(s))\|^{2}\right] \leq 2L_{J}^{2}\epsilon$$

And so

$$\sum_{j=0}^{k-1} \int_{j\epsilon}^{(j+1)\epsilon} \mathbb{E} \|\pi_{0,\gamma}(\bar{\alpha}(s))\nabla(J(\alpha(0)))\|^2$$
$$< 2k\epsilon L_I^2 \epsilon$$

Fourth Term: By Lemma 7.6, we have

$$\mathbb{E} \| K_{\Delta}(\theta_{\bar{s}}, \bar{\alpha}(s)) \hat{\nabla} J(\theta_{\bar{s}}) \|^2 \le 12\epsilon (L^2_{\nabla J}(M_{\theta} + 2B^2 M_{\theta}) + B^2 + \zeta) = 4\epsilon C_0$$

and so

(7.20)
$$\sum_{j=0}^{k-1} \int_{j\epsilon}^{(j+1)\epsilon} \mathbb{E} \| K_{\Delta}(\theta_{\bar{s}}, \bar{\alpha}(s)) \hat{\nabla} J(\theta_{\bar{s}}) \|^2 ds \le k\epsilon (4\epsilon C_0)$$

Combining (7.17), (7.18), (7.19), (7.20) in (7.14), we obtain

$$\begin{split} D(\mathbb{P}_{X}^{t} \| \mathbb{P}_{A}^{t}) \\ &\leq \beta L_{\nabla J}^{2} k \epsilon \left(\epsilon^{4} \left(72C_{0} + 18 \right) + \epsilon^{3} \left(6\sqrt{C_{0}} + \sqrt{2} \right) \right) \\ &+ 3\beta (k\epsilon)^{3} L_{\nabla J}^{2} \left(\epsilon^{4} \left(72C_{0} + 18 \right) + \epsilon^{3} \left(6\sqrt{C_{0}} + \sqrt{2} \right) \right) \\ &+ 2k\epsilon L_{J}^{2} \epsilon \\ &+ k\epsilon \left(4\epsilon C_{0} \right) \\ &\leq (k\epsilon)^{3} \epsilon^{3} \left[4\beta L_{\nabla J}^{2} \left(72C_{0} + 6\sqrt{C_{0}} + 18 + \sqrt{2} \right) \right] + (k\epsilon) \epsilon \left(2L_{J}^{2} + 4C_{0} \right) \end{split}$$

Now since $\pi_k = \text{Law}(\alpha_k)$ and $\nu_{k\epsilon} = \text{Law}(\alpha(t))$, the KL divergence data-processing inequality yields

$$D(\pi_k \| \nu_{k\epsilon}) \le D(\mathbb{P}_X^t \| \mathbb{P}_A^t) \qquad \Box$$

7.5. Proof of Proposition 5.3.

Proof. Recall that the continuous time diffusion of interest (2.12) has infinitesimal generator \mathcal{L} acting on C^2 function f as

$$\mathcal{L}f = \frac{1}{2}\pi_{0,\gamma}^2 \Delta f - \frac{\beta}{2}\pi_{0,\gamma} \langle \nabla J, \nabla f \rangle - \pi_{0,\gamma} \langle \nabla \pi_{0,\gamma}, \nabla f \rangle$$

We will show that the conditions of Proposition 4.4 hold:

1. Consider the Lyapunov function

$$V(x) = \exp\left(\frac{\beta m \|x\|^2}{2(\bar{\pi}^2_{0,\gamma} + 1)}\right)$$
36

Then we have

$$\begin{split} \mathcal{L}V(x) &= -\frac{\beta}{2} \pi_{0,\gamma}^2(x) \langle \nabla J(x), \nabla V(x) \rangle - \pi_{0,\gamma}(x) \langle \nabla \pi_{0,\gamma}(x), \nabla V(x) \rangle \\ &+ \frac{1}{2} \pi_{0,\gamma}^2(x) \Delta V(x) \\ &= \left\{ -\frac{\beta}{2} \frac{m\beta}{(\bar{\pi}_{0,\gamma}^2 + 1)} \pi_{0,\gamma}^2(x) \langle \nabla J(x), x \rangle - \frac{m\beta}{\bar{\pi}_{0,\gamma}^2 + 1} \pi_{0,\gamma}(x) \langle \nabla \pi_{0,\gamma}(x), x \rangle \right. \\ &+ \frac{1}{2} \pi_{0,\gamma}^2(x) (\frac{m\beta N}{\bar{\pi}_{0,\gamma}^2 + 1} + (\frac{m\beta}{\bar{\pi}_{0,\gamma}^2 + 1})^2 \|x\|^2)) \right\} V(x) \\ &\leq \left\{ -\frac{\beta}{2} \frac{m\beta}{\bar{\pi}_{0,\gamma}^2 + 1} \pi_{0,\gamma}^2(x) (m\|x\|^2 - b) + \frac{m\beta I}{\bar{\pi}_{0,\gamma}^2 + 1} \pi_{0,\gamma}(x) \\ &+ \frac{1}{2} \pi_{0,\gamma}^2(x) (\frac{m\beta N}{\bar{\pi}_{0,\gamma}^2 + 1} + (\frac{m\beta}{\bar{\pi}_{0,\gamma}^2 + 1})^2 \|x\|^2)) \right\} V(x) \\ &\leq \left\{ \left(\frac{1}{2} \beta mN + m\beta I + \frac{1}{2} \pi_{0,\gamma}^2(x) \frac{\beta^2 mb}{\bar{\pi}_{0,\gamma}^2 + 1} \right) \\ &- \left(\frac{\beta^2 m}{2(\bar{\pi}_{0,\gamma}^2 + 1)} \pi_{0,\gamma}^2(x) m - \frac{1}{2} \pi_{0,\gamma}^2(x) (\frac{\beta m}{\bar{\pi}_{0,\gamma}^2 + 1})^2 \right) \|x\|^2 \right\} V(x) \\ &\leq \left\{ \left(\frac{1}{2} \beta mN + m\beta I \right) + \frac{1}{2} \pi_{0,\gamma}^2(x) (\frac{\beta m}{\bar{\pi}_{0,\gamma}^2 + 1} \right) \\ &- \left(\frac{1}{2} \pi_{0,\gamma}^2(x) \frac{(\beta m)^2}{\bar{\pi}_{0,\gamma}^2 + 1} - \frac{1}{2} \pi_{0,\gamma}^2(x) (\frac{\beta m}{\bar{\pi}_{0,\gamma}^2 + 1})^2 \right) \|x\|^2 \right\} V(x) \\ &\leq \left\{ \left(\frac{1}{2} \beta mN + m\beta I \right) + \frac{1}{2} \pi_{0,\gamma}^2(x) \frac{\beta^2 mb}{\bar{\pi}_{0,\gamma}^2 + 1} \right) \\ &- \left(\frac{1}{2} \pi_{0,\gamma}^2(x) \frac{(\beta m)^2}{\bar{\pi}_{0,\gamma}^2 + 1} \left(1 - \frac{1}{\bar{\pi}_{0,\gamma}^2 + 1} \right) \right) \|x\|^2 \right\} V(x) \\ &\leq \left\{ \left(\frac{1}{2} \beta mN + m\beta I \right) + \left[\frac{1}{2} \pi_{0,\gamma}^2(x) \frac{\beta^2 mb}{\bar{\pi}_{0,\gamma}^2 + 1} \right) \\ &- \left(\frac{1}{2} \pi_{0,\gamma}^2(x) \left(\frac{(\beta m)^2}{\bar{\pi}_{0,\gamma}^2 + 1} M^2 \right) + 1 \\ &- \left(\frac{1}{2} \pi_{0,\gamma}^2(x) \left(\frac{(\beta m)^2}{\bar{\pi}_{0,\gamma}^2 + 1} M^2 \right) + 1 \\ &- \left(\frac{1}{2} \pi_{0,\gamma}^2(x) \left(\frac{(\beta m)^2}{\bar{\pi}_{0,\gamma}^2 + 1} \left(1 - \frac{1}{\bar{\pi}_{0,\gamma}^2 + 1} \right) \right) \|x\|^2 \right\} V(x) \end{aligned}$$

where in the last inequality statement we append

$$\frac{1}{2}\pi_{0,\gamma}^2(x)\left(\frac{(\beta m)^2}{\bar{\pi}_{0,\gamma}^2+1}M^2\right)+1$$

for the following reason: By assumption 3.4 we have that $\exists M : \pi_{0,\gamma}(x) < \frac{2}{(\beta m)^2 ||x||^2} \forall ||x|| > M, \gamma \leq 1$. We aim to show that the term inside brackets in the last inequality line of (1) is positive for all x. First take ||x|| < M: we 37

have that:

$$\left(\frac{1}{2}\pi_{0,\gamma}^2(x)\frac{(\beta m)^2}{\bar{\pi}_{0,\gamma}^2+1}\left(1-\frac{1}{\bar{\pi}_{0,\gamma}^2+1}\right)\right)\|x\|^2 < \frac{1}{2}\pi_{0,\gamma}^2(x)\left(\frac{(\beta m)^2}{\bar{\pi}_{0,\gamma}^2+1}M^2\right)$$

Now consider $||x|| \ge M$. By Lemma 7.9 we have

$$\begin{split} &\frac{1}{2}\pi_{0,\gamma}(x)(\beta m)^2 \|x\|^2 < 1\\ &\Rightarrow \frac{1}{2}\pi_{0,\gamma}^2(x)\frac{(\beta m)^2}{\bar{\pi}_{0,\gamma}^2 + 1}\|x\|^2 < 1\\ &\Rightarrow \frac{1}{2}\pi_{0,\gamma}^2(x)\frac{(\beta m)^2}{\bar{\pi}_{0,\gamma}^2 + 1}\left(1 - \frac{1}{\bar{\pi}_{0,\gamma}^2 + 1}\right)\|x\|^2 < 1 \end{split}$$

Thus we have:

$$\begin{split} & \left[\frac{1}{2}\pi_{0,\gamma}^{2}(x)\frac{\beta^{2}mb}{\bar{\pi}_{0,\gamma}^{2}+1} + \frac{1}{2}\pi_{0,\gamma}^{2}(x)\left(\frac{(\beta m)^{2}}{\bar{\pi}_{0,\gamma}^{2}+1}M^{2}\right) + 1 \\ & -\left(\frac{1}{2}\pi_{0,\gamma}^{2}(x)\frac{(\beta m)^{2}}{\bar{\pi}_{0,\gamma}^{2}+1}\left(1 - \frac{1}{\bar{\pi}_{0,\gamma}^{2}+1}\right)\right)\|x\|^{2}\right] > 0 \quad \forall x \in \mathbb{R}^{N} \end{split}$$

and now observe that

$$\begin{split} & \left[\frac{1}{2}\pi_{0,\gamma}^{2}(x)\frac{\beta^{2}mb}{\bar{\pi}_{0,\gamma}^{2}+1} + \frac{1}{2}\pi_{0,\gamma}^{2}(x)\left(\frac{(\beta m)^{2}}{\bar{\pi}_{0,\gamma}^{2}+1}M^{2}\right) + 1 \\ & - \left(\frac{1}{2}\pi_{0,\gamma}^{2}(x)\frac{(\beta m)^{2}}{\bar{\pi}_{0,\gamma}^{2}+1}\left(1 - \frac{1}{\bar{\pi}_{0,\gamma}^{2}+1}\right)\right)\|x\|^{2}\right] \\ & \leq \left[\frac{1}{2}(\pi_{0,\gamma}^{2}(x)+1)\frac{\beta^{2}mb}{\bar{\pi}_{0,\gamma}^{2}+1} + \frac{1}{2}(\pi_{0,\gamma}^{2}(x)+1)\left(\frac{(\beta m)^{2}}{\bar{\pi}_{0,\gamma}^{2}+1}M^{2}\right) \right. \\ & - \left(\frac{1}{2}(\pi_{0,\gamma}^{2}(x)+1)\left[\frac{1}{2}\frac{\beta^{2}mb}{\bar{\pi}_{0,\gamma}^{2}+1}\left(1 - \frac{1}{\bar{\pi}_{0,\gamma}^{2}+1}M^{2}\right)\right] \\ & \leq (\pi_{0,\gamma}^{2}(x)+1)\left[\frac{1}{2}\frac{\beta^{2}mb}{\bar{\pi}_{0,\gamma}^{2}+1} + \frac{1}{2}\left(\frac{(\beta m)^{2}}{\bar{\pi}_{0,\gamma}^{2}+1}M^{2}\right) \\ & - \left(\frac{1}{2}\frac{(\beta m)^{2}}{\bar{\pi}_{0,\gamma}^{2}+1}\left(1 - \frac{1}{\bar{\pi}_{0,\gamma}^{2}+1}\right)\right)\|x\|^{2}\right] \\ & \leq (\bar{\pi}_{0,\gamma}^{2}+1)\left[\frac{1}{2}\frac{\beta^{2}mb}{\bar{\pi}_{0,\gamma}^{2}+1} + \frac{1}{2}\left(\frac{(\beta m)^{2}}{\bar{\pi}_{0,\gamma}^{2}+1}M^{2}\right) \\ & - \left(\frac{1}{2}\frac{(\beta m)^{2}}{\bar{\pi}_{0,\gamma}^{2}+1}\left(1 - \frac{1}{\bar{\pi}_{0,\gamma}^{2}+1}\right)\right)\|x\|^{2}\right] \end{split}$$

where we use that $1 \leq \frac{(\beta m)^2}{\bar{\pi}_{0,\gamma+1}}M^2$, which is derived from Assumption (3.8). Thus we have:

$$\frac{\mathcal{L}V(x)}{V(x)} \le \left\{ \left(\frac{1}{2}\beta mN + m\beta I\right) \right\}$$
38

$$\begin{split} &+ (\bar{\pi}_{0,\gamma}^2 + 1) \left[\frac{1}{2} \frac{\beta^2 m b}{\bar{\pi}_{0,\gamma}^2 + 1} + \frac{1}{2} \left(\frac{(\beta m)^2}{\bar{\pi}_{0,\gamma}^2 + 1} M^2 \right) \right] \\ &- (\bar{\pi}_{0,\gamma}^2 + 1) \left(\frac{1}{2} \frac{(\beta m)^2}{\bar{\pi}_{0,\gamma}^2 + 1} \left(1 - \frac{1}{\bar{\pi}_{0,\gamma}^2 + 1} \right) \right) \|x\|^2 \right\} \\ &\leq \left(\frac{1}{2} \beta m N + \beta m I \right) + \frac{1}{2} \left[\beta^2 m b + (\beta m M)^2 \right] \\ &- \frac{1}{2} \left((\beta m)^2 + \left(1 - \frac{1}{i dist max^2 + 1} \right) \right) \|x\|^2 \\ &:= \kappa - \gamma \|x\|^2 \end{split}$$

and observe that $\kappa, \gamma > 0$. Thus, condition (1) of Proposition 4.4 holds. 2. From (1), we have

$$\frac{\mathcal{L}V(x)}{V(x)} \le \kappa - \gamma \|x\|^2$$

Observe that this implies

$$\frac{\mathcal{L}V(x)}{V(x)} \le -\kappa + 2\kappa \mathbf{1}_{(\|x\|^2 \le 2\kappa/\gamma)}$$

Moreover, by Lemma 7.10 and since $J(x) \ge 0$, we have

$$O_r(\beta J) \le \beta \left(\frac{L_{\nabla J}}{2} \|x\|^2 + \|x\| + A\right) \le \beta \left(\frac{(L_{\nabla J} + B)r^2}{2} + A + B\right)$$

and so by Proposition 4.3, with $\kappa_0 = 2\kappa$, $\zeta_0 = \kappa$, $r^2 = 2\kappa/\gamma$, we have that π_∞ satisfies a Poincarê inequality with constant

$$c_P = \frac{1}{\lambda} \le \frac{1}{2\kappa} \left(1 + \frac{4C\kappa^2}{\gamma} \exp\left(\beta \left(\frac{(L_{\nabla J} + B)\kappa}{\gamma} + A + B\right)\right) \right)$$

where κ and γ are defined above and provided in (5.6).

3. By assumption 3.1, we have

$$\nabla^2 \beta J \succeq -\beta L_{\nabla J} I_d \succeq 0$$

Thus the conditions of Proposition 4.4 are met, with K = 0. So, letting $\zeta = 1$:

$$Z_1 = \frac{2}{\gamma} + 1, \quad Z_2 = \frac{2}{\gamma} \left(\kappa + \gamma \int_{\mathbb{R}^N} \|w\|^2 \pi_\infty(dw) \right)$$

Now we would like to make the bound on c_{LS} (5.5) more explicit by providing a bound on $\int_{\mathbb{R}^N} ||w||^2 \pi_{\infty}(dw)$. From (5.8) we have $\mathcal{W}_2(\nu_t, \pi_{\infty}) \to 0$ as $t \to \infty$, and thus by Theorem 7.12 of [27] and Lemma 7.5 it follows that (with $\gamma \leq 1$)

(7.21)
$$\int_{\mathbb{R}^{N}} \|w\|^{2} \pi_{\infty}(dw) = \lim_{t \to \infty} \int_{\mathbb{R}^{N}} \|w\|^{2} \nu_{t}(dw)$$
$$\leq \kappa_{0}^{\gamma} + \frac{(\beta b + N)\bar{\pi}_{0,\gamma} + 2I}{(m\beta)\bar{\pi}_{0,\gamma}} \leq \kappa_{0} + \frac{(\beta b + N)\bar{\pi} + 2I}{(m\beta)\bar{\pi}}$$

So now letting

$$Z_2 = \frac{2}{\gamma} \left(\kappa + \gamma \left(\kappa_0 + \frac{(\beta b + N)\bar{\pi} + 2I}{(m\beta)\bar{\pi}} \right) \right)$$

we have that π_∞ satisfies a log-Sobolev inequality with constant

$$c_{LS} = Z_1 + (Z_2 + 2)c_P$$