

Stochastic Gradient Descent in Continuous Time

Justin Sirignano* and Konstantinos Spiliopoulos^{†‡}

November 20, 2016

Abstract

We consider stochastic gradient descent for continuous-time models. Traditional approaches for the statistical estimation of continuous-time models, such as batch optimization, can be impractical for large datasets where observations occur over a long period of time. Stochastic gradient descent provides a computationally efficient method for such statistical estimation problems. The stochastic gradient descent algorithm performs an online parameter update in continuous time, with the parameter updates satisfying a stochastic differential equation. The parameters are proven to converge to a local minimum of a natural objective function for the estimation of the continuous-time dynamics. The convergence proof leverages ergodicity by using an appropriate Poisson equation to help describe the evolution of the parameters for large times. Numerical analysis of the stochastic gradient descent algorithm is presented for several examples, including the Ornstein-Uhlenbeck process, Burger's stochastic partial differential equation, and reinforcement learning.

1 Introduction

Batch optimization for the statistical estimation of continuous-time models can be impractical for large datasets where observations occur over a long period of time. Batch optimization takes a sequence of descent steps for the model error for the entire observed data path. Since each descent step is for the model error for the *entire observed data path*, batch optimization is slow (sometimes impractically slow) for long periods of time or models which are computationally costly to evaluate (e.g., partial differential equations).

Stochastic gradient descent in continuous time provides a computationally efficient method for statistical learning over long time periods and for complex models. Stochastic gradient descent *continually* takes gradient steps *along the path of the observation* which results in much more rapid convergence. Parameters are updated online in continuous time, with the parameter updates satisfying a stochastic differential equation. We prove that the parameters converge to a local minimum of a natural objective function for the estimation of the continuous-time dynamics.

We consider a diffusion $X_t \in \mathcal{X} = \mathbb{R}^m$:

$$dX_t = f^*(X_t)dt + \sigma dW_t. \quad (1.1)$$

The goal is to statistically estimate a model $f(x, \theta)$ for $f^*(x)$ where $\theta \in \mathbb{R}^n$. W_t is a standard Brownian motion. The diffusion term W_t represents any random behavior of the system or environment. The functions $f(x, \theta)$ and $f^*(x)$ may be non-convex.

The stochastic gradient descent update in continuous time follows the stochastic differential equation (SDE):

$$d\theta_t = \alpha_t [\nabla_{\theta} f(X_t; \theta_t)(\sigma\sigma^T)^{-1}dX_t - \nabla_{\theta} f(X_t, \theta_t)(\sigma\sigma^T)^{-1}f(X_t, \theta_t)dt], \quad (1.2)$$

where $\nabla_{\theta} f(X_t; \theta_t)$ is matrix valued and α_t is the learning rate. The parameter update (1.2) can be used for both statistical estimation given previously observed data as well as online learning (i.e., statistical estimation in real-time as data becomes available).

*Department of Industrial and Enterprise Systems Engineering, University of Illinois at Urbana Champaign, Urbana, Email: jasilrign@illinois.edu

[†]Department of Mathematics and Statistics, Boston University, Boston, E-mail: kspiliop@math.bu.edu

[‡]Research of K.S. supported in part by the National Science Foundation (DMS 1550918)

We assume that X_t is sufficiently ergodic (to be concretely specified later in the paper) and that it has some well-behaved $\pi(dx)$ as its unique invariant measure. As a general notation, if $h(x, \theta)$ is a generic $L^1(\pi)$ function, then we define its average over $\pi(dx)$ to be

$$\bar{h}(\theta) = \int_{\mathcal{X}} h(x, \theta) \pi(dx)$$

Let us set

$$g(x, \theta) = \frac{1}{2} \|f(x, \theta) - f^*(x)\|_{\sigma\sigma^T}^2 = \frac{1}{2} \left\langle f(x, \theta) - f^*(x), (\sigma\sigma^T)^{-1} (f(x, \theta) - f^*(x)) \right\rangle$$

Heuristically, it is expected that θ_t will tend towards the minimum of the function $\bar{g}(\theta) = \int_{\mathcal{X}} g(x, \theta) \pi(dx)$. The stochastic gradient descent update (1.2) continuously moves the parameters in an *estimated* direction of $\nabla_{\theta} \bar{g}(\theta)$. That is, the drift of θ_t is a *biased* estimate of $\nabla_{\theta} \bar{g}(\theta)$. The bias will decrease as time increases. In the standard discrete case of stochastic gradient descent where data is i.i.d. at every step, the gradient updates are unbiased. The data X_t in the continuous time will be correlated over long periods of time, further complicating the analysis.

In this paper we show that if α_t is appropriately chosen then $\nabla_{\theta} \bar{g}(\theta_t) \rightarrow 0$ as $t \rightarrow \infty$, see Theorem 2.4. Results like this have been previously derived in the literature for discrete time and in the absence of the X component, see [1]. The presence of the X term complicates the analysis as one needs to control the speed at which convergence to equilibrium happens. Furthermore, the parameter updates are now biased and the bias is correlated across times.

Although stochastic gradient descent for discrete time has been extensively studied, stochastic gradient descent in continuous time has received relatively little attention. In comparison to results available in discrete time, our convergence result requires weaker assumptions. We refer readers to [6] and [1] for a thorough review of the very large literature on stochastic gradient descent. There are also many algorithms which modify traditional stochastic gradient descent (stochastic gradient descent with momentum, Adagrad, RMSprop, etc.). For a review of these variants of stochastic gradient descent, see [8]. We mention below the prior work which is most relevant to our paper.

As mentioned, [1] proves convergence of stochastic gradient descent in discrete time in the absence of the X process. The presence of the X process is essential for considering a wide range of problems in continuous time, and showing convergence with its presence is considerably more difficult. The X term introduces correlation across times, and this correlation does not disappear as time tends to infinity. Unlike in [1] where parameter updates are unbiased, the correlation introduced by the X process causes parameter updates to be biased. Furthermore, the bias of the parameter updates is correlated across times. These facts make it challenging to prove convergence in the continuous-time case. In order to prove convergence, we use an appropriate Poisson equation associated with X to describe the evolution of the parameters for large times.

Notably, [9] proves convergence in L^2 of projected stochastic gradient descent for convex functions (in a set-up different to ours). In projected gradient descent, the parameters are projected back into an a priori chosen compact set. Therefore, the algorithm cannot hope to reach the minimum if the minimum is located outside of the chosen compact set. Of course, the compact set can be chosen to be very large for practical purposes. Our paper considers unconstrained stochastic gradient descent of not necessarily convex functions and proves almost sure convergence taking into account the X component as well.

Another approach for proving convergence is to show that the algorithm converges to the solution of an ODE which itself converges to a limiting point; see [6] and [7]. This method, sometimes called the “ODE method”, requires the strong assumption that the iterates (i.e., the model parameters which are being learned) remain in a bounded set with probability one. Proving that the iterates remain in a bounded set with probability one can be challenging to show and, moreover, may not necessarily be true for many models of interest. It is also noteworthy that it is not clear that the ODE method can be used to prove convergence of the stochastic gradient descent algorithm in continuous time, which is the problem that this paper considers.

The paper [10] studies continuous-time stochastic mirror descent in a setting different than ours. They prove a bound for the minimization of a convex function. In the framework of [10], the objective function is known and descent steps are therefore unbiased. In contrast, we consider the statistical estimation of the

unknown dynamics of a random process (i.e. the X process satisfying (1.1)). The descent steps in our case are therefore biased, complicating the analysis.

For certain applications, stochastic gradient descent in continuous time may have advantages over stochastic gradient descent in discrete time. Physics and engineering models are typically in continuous time. It therefore makes sense to also develop the statistical learning updates in continuous time. For example, statistical learning may be used to estimate coefficients or parameters in the engineering model. The engineering model may also be enhanced in some cases with the addition of a machine learning model to better fit real-world conditions and data.

Continuous-time dynamics are oftentimes simpler to analyze than discrete dynamics at longer time intervals. For instance, a partial differential equation can be written in a very simple form, but its global dynamics over a long time period can be very complex. It may therefore be advantageous to learn the continuous-time dynamics.

In addition, the continuous-time framework only requires the estimation of a deterministic function $f(x, \theta)$ while a discrete-time framework would require estimating a multidimensional density $p(x, x', \theta) = \mathbb{P}[X_{t+\Delta} \in dx | X_t = x', \theta]$. Estimating a multidimensional density $p(x, x', \theta)$ is typically more computationally challenging than estimating a vector-valued function $f(x, \theta)$. In our setting (1.1), the instantaneous dynamics $f(x, \theta)$ directly yield the density $p(x, x', \theta)$. Consequently, even if one's goal is to model the density at a longer time horizon $\Delta \gg 0$, it may be preferable to instead estimate the continuous-time dynamics f instead of directly estimating the density p .

Although continuous-time stochastic gradient descent must ultimately be discretized for numerical implementation, there are still significant numerical advantages to the continuous-time formulation. The continuous-time stochastic gradient descent algorithm allows for the control and reduction of numerical error due to discretization. In particular, higher-order numerical schemes can be used for the numerical solution of the continuous-time stochastic gradient descent updates. This will lead to more accurate and more computationally efficient parameter updates. Furthermore, continuous-time stochastic gradient descent can use non-uniform time step sizes. If convergence is slow, the time step size may be adaptively decreased. Continuous-time stochastic gradient descent allows for efficient learning at different time step sizes by providing the appropriate scaling for the learning rate.¹ In contrast, discrete-time stochastic gradient descent operates at fixed discrete steps and does not allow for the adaptive control of the time step size nor higher-order numerical schemes to reduce discretization error.

We numerically study the convergence of continuous stochastic gradient descent for a number of applications. Applications include the Ornstein-Uhlenbeck process, Burger's equation, and the classic reinforcement learning problem of balancing a pole on a moving cart. The Ornstein-Uhlenbeck process is widely-used in finance, physics, and biology. Burger's equation is a widely used nonlinear partial differential equation which is important to fluid mechanics, nonlinear acoustics, and aerodynamics. It is extensively used in engineering.

The paper is organized into three main sections. Section 2 presents the assumption and the main theorem. In Section 3 we prove the main result of this paper on the convergence of continuous stochastic gradient descent. Section 4 provides numerical analysis of continuous stochastic gradient descent for several example applications.

2 Assumptions and Main Result

Before presenting the main result of this paper, Theorem 2.4, let us elaborate on the standing assumptions. In regards to the learning rate α_t the standing assumption is

Condition 2.1. Assume that $\int_0^\infty \alpha_t dt = \infty$, $\int_0^\infty \alpha_t^2 dt < \infty$ and that $\int_0^\infty |\alpha'_s| ds < \infty$.

A standard choice for α_t that satisfies Condition 2.1 is $\alpha_t = \frac{1}{C+t}$ for some constant $0 < C < \infty$. Notice that the condition $\int_0^\infty |\alpha'_s| ds < \infty$ follows immediately from the other two restrictions for the learning rate if it is chosen to be a monotonic function of t .

Let us next discuss the assumptions that we impose on σ , $f^*(x)$ and $f(x, \theta)$. Condition 2.2 guarantees uniqueness and existence of an invariant measure for the X process.

Condition 2.2. We assume that σ is non-degenerate bounded diffusion matrix and $\lim_{|x| \rightarrow \infty} f^*(x) \cdot x = -\infty$

¹For instance, the learning rate after an Euler discretization is $\alpha_t \Delta t$ where Δt is the time step size.

In addition, with respect to $\nabla_\theta f(x, \theta)$ we assume that $\theta \in \mathbb{R}^n$ and we impose the following condition

Condition 2.3. 1. We assume that $\nabla_\theta g(x, \cdot) \in C^2(\mathbb{R}^n)$ for all $x \in \mathcal{X}$, $\frac{\partial^2 \nabla_\theta g}{\partial x^2} \in C(\mathcal{X}, \mathbb{R}^n)$, $\nabla_\theta g(\cdot, \theta) \in C^\alpha(\mathcal{X})$ uniformly in $\theta \in \mathbb{R}^n$ for some $\alpha \in (0, 1)$ and that there exist K and q such that

$$\sum_{i=0}^2 \left| \frac{\partial^i \nabla_\theta g}{\partial \theta^i}(x, \theta) \right| \leq K(1 + |x|^q).$$

2. For every $N > 0$ there exists a constant $C(N)$ such that for all $\theta_1, \theta_2 \in \mathbb{R}^n$ and $|x| \leq N$, the diffusion coefficient $\nabla_\theta f$ satisfies

$$|\nabla_\theta f(x, \theta_1) - \nabla_\theta f(x, \theta_2)| \leq C(N)|\theta_1 - \theta_2|.$$

Moreover, there exists $K > 0$ and $q > 0$ such that

$$|\nabla_\theta f(x, \theta)| \leq K(1 + |x|^q).$$

3. The function $f^*(x)$ is $C_b^{2+\alpha}(\mathcal{X})$ with $\alpha \in (0, 1)$. Namely, it has two bounded derivatives in x , with all partial derivatives being Hölder continuous, with exponent α , with respect to x .

In a sense, Condition 2.3, allow us to control the ergodic behavior of the X process. As it will be seen from the proof of the main convergence result, Theorem 2.4, one needs to control terms of the form $\int_0^t \alpha_s (\nabla \bar{g}(\theta_s) - g(X_s, \theta_s)) ds$. Due to ergodicity of the X process one expects that such terms are small in magnitude and go to zero as $t \rightarrow \infty$. However, the speed at which they go to zero is what matters here. We treat such terms by rewriting them equivalently using appropriate Poisson type partial differential equations (PDE). Then Condition 2.3, guarantees that such equations have unique solutions that do not grow faster than polynomially in the x variable, this is Theorem A.1 in Appendix A.

The main result of this paper is Theorem 2.4.

Theorem 2.4. *Assume that Conditions 2.1, 2.2 and 2.3 hold. Then we have that*

$$\lim_{t \rightarrow \infty} \|\nabla \bar{g}(\theta_t)\| = 0, \text{ almost surely.}$$

3 Proof of Theorem 2.4

We proceed in a spirit similar to that of [1]. However, apart from continuous versus discrete dynamics, one of the main challenges of the proof here is the presence of the ergodic X process. Let us consider an arbitrarily given $\kappa > 0$ and $\lambda = \lambda(\kappa) > 0$ to be chosen. Then set $\sigma_0 = 0$ and consider the cycles of random times

$$0 = \sigma_0 \leq \tau_1 \leq \sigma_1 \leq \tau_2 \leq \sigma_2 \leq \dots$$

where for $k = 1, 2, \dots$

$$\begin{aligned} \tau_k &= \inf\{t > \sigma_{k-1} : \|\nabla \bar{g}(\theta_t)\| \geq \kappa\}, \\ \sigma_k &= \sup\{t > \tau_k : \frac{\|\nabla \bar{g}(\theta_{\tau_k})\|}{2} \leq \|\nabla \bar{g}(\theta_s)\| \leq 2\|\nabla \bar{g}(\theta_{\tau_k})\| \text{ for all } s \in [\tau_k, t] \text{ and } \int_{\tau_k}^t \alpha_s ds \leq \lambda\}. \end{aligned}$$

The purpose of these random times is to control the periods of time where $\|\nabla \bar{g}(\theta)\|$ is close to zero and away from zero. Let us next define the random time intervals $J_k = [\sigma_{k-1}, \tau_k)$ and $I_k = [\tau_k, \sigma_k)$. Notice that for every $t \in J_k$ we have $\|\nabla \bar{g}(\theta_t)\| < \kappa$.

Let us next consider some $\eta > 0$ sufficiently small to be chosen later on and set $\sigma_{k,\eta} = \sigma_k + \eta$.

Lemma 3.1. *Assume that Conditions 2.1, 2.2 and 2.3 hold. Choose $\lambda > 0$ such that for a given $\kappa > 0$, one has $3\lambda + \frac{\lambda}{4\kappa} = \frac{1}{2L_{\nabla \bar{g}}}$, where $L_{\nabla \bar{g}}$ is the Lipschitz constant of $\nabla \bar{g}$. For k large enough and for $\eta > 0$ small enough (potentially random depending on k), one has $\int_{\tau_k}^{\sigma_{k,\eta}} \alpha_s ds > \lambda$. In addition we also have $\frac{\lambda}{2} \leq \int_{\tau_k}^{\sigma_k} \alpha_s ds \leq \lambda$ with probability one.*

Proof. Let us define the random variable

$$R_s = \sum_{k \geq 1} \|\nabla \bar{g}(\theta_{\tau_k})\| 1_{s \in I_k} + \kappa 1_{s \in [0, \infty) \setminus \bigcup_{k \geq 1} I_k}. \quad (3.1)$$

Then, for any $s \in \mathbb{R}$ we have $\|\nabla \bar{g}(\theta_s)\|/R_s \leq 2$.

We proceed with an argument via contradiction. In particular let us assume that $\int_{\tau_k}^{\sigma_{k,\eta}} \alpha_s ds \leq \lambda$ and let us choose arbitrarily some $\epsilon > 0$ such that $\epsilon \leq \lambda/8$.

Let us now make some remarks that are independent of the sign of $\int_{\tau_k}^{\sigma_{k,\eta}} \alpha_s ds - \lambda$. Due to the summability condition $\int_0^\infty \alpha_t^2 dt < \infty$, $\frac{\kappa}{\|\nabla \bar{g}(\theta_{\tau_k})\|} \leq 1$ and Conditions 2.1 and 2.3, the martingale convergence theorem applies to the martingale $\int_0^t \alpha_s \frac{\kappa}{\|\nabla \bar{g}(\theta_{\tau_k})\|} \nabla_\theta f(X_s, \theta_s) \sigma^{-1} dW_s$. This means that there exists a square integrable random variable M such that $\int_0^t \alpha_s \frac{\kappa}{\|\nabla \bar{g}(\theta_{\tau_k})\|} \nabla_\theta f(X_s, \theta_s) \sigma^{-1} dW_s \rightarrow M$ both almost surely and in L^2 . This means that for the given $\epsilon > 0$ there is k large enough such that $\left\| \int_{\tau_k}^{\sigma_{k,\eta}} \alpha_s \frac{\kappa}{\|\nabla \bar{g}(\theta_{\tau_k})\|} \nabla_\theta f(X_s, \theta_s) \sigma^{-1} dW_s \right\| < \epsilon$ almost surely.

Let us also assume that for the given k, η is so small such that for any $s \in [\tau_k, \sigma_{k,\eta}]$ one has $\|\nabla \bar{g}(\theta_s)\| \leq 3\|\nabla \bar{g}(\theta_{\tau_k})\|$.

Then, we obtain the following

$$\begin{aligned} \|\theta_{\sigma_{k,\eta}} - \theta_{\tau_k}\| &= \left\| - \int_{\tau_k}^{\sigma_{k,\eta}} \alpha_s \nabla_\theta g(X_s, \theta_s) ds + \int_{\tau_k}^{\sigma_{k,\eta}} \alpha_s \nabla_\theta f(X_s, \theta_s) \sigma^{-1} dW_s \right\| \\ &= \left\| - \int_{\tau_k}^{\sigma_{k,\eta}} \alpha_s \nabla_\theta \bar{g}(\theta_s) ds - \int_{\tau_k}^{\sigma_{k,\eta}} \alpha_s (\nabla_\theta g(X_s, \theta_s) - \nabla_\theta \bar{g}(\theta_s)) ds + \int_{\tau_k}^{\sigma_{k,\eta}} \alpha_s \nabla_\theta f(X_s, \theta_s) \sigma^{-1} dW_s \right\| \\ &\leq \int_{\tau_k}^{\sigma_{k,\eta}} \alpha_s \|\nabla \bar{g}(\theta_s)\| ds + \left\| \int_{\tau_k}^{\sigma_{k,\eta}} \alpha_s (\nabla_\theta g(X_s, \theta_s) - \nabla_\theta \bar{g}(\theta_s)) ds \right\| + \left\| \int_{\tau_k}^{\sigma_{k,\eta}} \alpha_s \nabla_\theta f(X_s, \theta_s) \sigma^{-1} dW_s \right\| \\ &\leq 3\|\nabla \bar{g}(\theta_{\tau_k})\| \int_{\tau_k}^{\sigma_{k,\eta}} \alpha_s ds + \left\| \int_{\tau_k}^{\sigma_{k,\eta}} \alpha_s (\nabla_\theta g(X_s, \theta_s) - \nabla_\theta \bar{g}(\theta_s)) ds \right\| \\ &\quad + \frac{\|\nabla \bar{g}(\theta_{\tau_k})\|}{\kappa} \left\| \int_{\tau_k}^{\sigma_{k,\eta}} \alpha_s \frac{\kappa}{\|\nabla \bar{g}(\theta_{\tau_k})\|} \nabla_\theta f(X_s, \theta_s) \sigma^{-1} dW_s \right\| \\ &\leq 3\|\nabla \bar{g}(\theta_{\tau_k})\| \lambda + \left\| \int_{\tau_k}^{\sigma_{k,\eta}} \alpha_s (\nabla_\theta g(X_s, \theta_s) - \nabla_\theta \bar{g}(\theta_s)) ds \right\| + \frac{\|\nabla \bar{g}(\theta_{\tau_k})\|}{\kappa} \epsilon \\ &\leq 3\|\nabla \bar{g}(\theta_{\tau_k})\| \lambda + \left\| \int_{\tau_k}^{\sigma_{k,\eta}} \alpha_s (\nabla_\theta g(X_s, \theta_s) - \nabla_\theta \bar{g}(\theta_s)) ds \right\| + \frac{\|\nabla \bar{g}(\theta_{\tau_k})\|}{\kappa} \lambda/8 \\ &= \|\nabla \bar{g}(\theta_{\tau_k})\| \left[3\lambda + \frac{\lambda}{8\kappa} \right] + \left\| \int_{\tau_k}^{\sigma_{k,\eta}} \alpha_s (\nabla_\theta g(X_s, \theta_s) - \nabla_\theta \bar{g}(\theta_s)) ds \right\|. \end{aligned}$$

Let us next bound appropriately the Euclidean norm of the vector-valued random variable

$$\Gamma_{k,\eta} = \int_{\tau_k}^{\sigma_{k,\eta}} \alpha_s (\nabla_\theta g(X_s, \theta_s) - \nabla_\theta \bar{g}(\theta_s)) ds.$$

By Lemma 3.2 we have that for the same $0 < \epsilon < \lambda/8$ that was chosen before there is k large enough such that almost surely

$$\|\Gamma_{k,\eta}\| \leq \epsilon \leq \lambda/8.$$

Hence, using also the fact that $\frac{\kappa}{\|\nabla \bar{g}(\theta_{\tau_k})\|} \leq 1$ we obtain

$$\|\theta_{\sigma_{k,\eta}} - \theta_{\tau_k}\| \leq \|\nabla \bar{g}(\theta_{\tau_k})\| \left[3\lambda + \frac{\lambda}{4\kappa} \right] = \|\nabla \bar{g}(\theta_{\tau_k})\| \frac{1}{2L_{\nabla \bar{g}}}.$$

The latter then implies that we should have

$$\|\nabla \bar{g}(\theta_{\sigma_{k,\eta}}) - \nabla \bar{g}(\theta_{\tau_k})\| \leq L_{\nabla \bar{g}} \|\theta_{\sigma_{k,\eta}} - \theta_{\tau_k}\| \leq \frac{\|\nabla \bar{g}(\theta_{\tau_k})\|}{2}.$$

The latter statement will then imply that

$$\frac{\|\nabla \bar{g}(\theta_{\tau_k})\|}{2} \leq \|\nabla \bar{g}(\theta_{\sigma_{k,\eta}})\| \leq 2\|\nabla \bar{g}(\theta_{\tau_k})\|.$$

But then we would necessarily have that $\int_{\tau_k}^{\sigma_{k,\eta}} \alpha_s ds > \lambda$, since otherwise $\sigma_{k,\eta} \in [\tau_k, \sigma_k]$ which is impossible.

Next we move on to prove the second statement of the lemma. By definition we have $\int_{\tau_k}^{\sigma_k} \alpha_s ds \leq \lambda$. So it remains to show that $\frac{\lambda}{2} \leq \int_{\tau_k}^{\sigma_k} \alpha_s ds$. Since we know that $\int_{\tau_k}^{\sigma_{k,\eta}} \alpha_s ds > \lambda$ and because for k large enough and η small enough one should have $\int_{\sigma_k}^{\sigma_{k,\eta}} \alpha_s ds \leq \lambda/2$, we obtain that

$$\int_{\tau_k}^{\sigma_k} \alpha_s ds \geq \lambda - \int_{\sigma_k}^{\sigma_{k,\eta}} \alpha_s ds \geq \lambda - \lambda/2 = \lambda/2,$$

concluding the proof of the lemma. \square

Lemma 3.2. *Assume that Conditions 2.1, 2.2 and 2.3 hold. Let us set*

$$\Gamma_{k,\eta} = \int_{\tau_k}^{\sigma_{k,\eta}} \alpha_s (\nabla_\theta g(X_s, \theta_s) - \nabla_\theta \bar{g}(\theta_s)) ds.$$

Then, with probability one we have that

$$\|\Gamma_{k,\eta}\| \rightarrow 0, \text{ as } k \rightarrow \infty.$$

Proof. The idea is to use Theorem A.1 in order to get an equivalent expression for the term $\Gamma_{k,\eta}$ that we seek to control.

Let us consider the function $G(t, x, \theta) = \alpha_t (\nabla_\theta g(x, \theta) - \nabla_\theta \bar{g}(\theta))$. Notice that by definition and due to Condition 2.3, the function $G(t, x, \theta)$ satisfies the centering condition (A.1) of Theorem A.1 componentwise. So, the Poisson equation (A.2) will have a unique smooth solution that grows at most polynomially in x . Let us apply Itô formula to the vector valued function $u(t, x, \theta)$ that is solution to this Poisson equation with right hand side $G(t, x, \theta)$. Doing so, we get for $i = 1, \dots, n$

$$\begin{aligned} u_i(\sigma, X_\sigma, \theta_\sigma) - u_i(\tau, X_\tau, \theta_\tau) &= \int_\tau^\sigma \partial_s u_i(s, X_s, \theta_s) ds + \int_\tau^\sigma \mathcal{L}_x u_i(s, X_s, \theta_s) ds + \int_\tau^\sigma \mathcal{L}_\theta u_i(s, X_s, \theta_s) ds \\ &\quad + \int_\tau^\sigma \alpha_s \text{tr} [\nabla_\theta f(X_s, \theta_s) \nabla_x \nabla_\theta u_i(s, X_s, \theta_s)] ds \\ &\quad + \int_\tau^\sigma \langle \nabla_x u_i(s, X_s, \theta_s), \sigma dW_s \rangle + \int_\tau^\sigma \alpha_s \langle \nabla_\theta u_i(s, X_s, \theta_s), \nabla_\theta f(X_s, \theta_s) \sigma^{-1} dW_s \rangle, \end{aligned}$$

where \mathcal{L}_x and \mathcal{L}_θ denote the infinitesimal generators for processes X and θ respectively.

Recall now that $u(s, x, \theta)$ is solution to the given Poisson equation and that we can write $u(s, x, \theta) = \alpha_s v(x, \theta)$. Using these facts and rearranging the previous Itô formula, we get in vector notation

$$\begin{aligned} \Gamma_{k,\eta} &= \int_{\tau_k}^{\sigma_k} \alpha_s (\nabla_\theta g(X_s, \theta_s) - \nabla_\theta \bar{g}(\theta_s)) ds = \int_{\tau_k}^{\sigma_k} \mathcal{L}_x u(s, X_s, \theta_s) ds \\ &= \left[\alpha_{\sigma_k} v(X_{\sigma_k}, \theta_{\sigma_k}) - \alpha_{\tau_k} v(X_{\tau_k}, \theta_{\tau_k}) - \int_{\tau_k}^{\sigma_k} \partial_s \alpha_s v(X_s, \theta_s) ds \right] \\ &\quad - \int_{\tau_k}^{\sigma_k} \alpha_s [\mathcal{L}_\theta v(X_s, \theta_s) + \alpha_s \text{tr} [\nabla_\theta f(X_s, \theta_s) \nabla_{x_i} \nabla_\theta v(X_s, \theta_s)]_{i=1}^m] ds \\ &\quad - \int_{\tau_k}^{\sigma_k} \alpha_s \langle \nabla_x v(X_s, \theta_s), \sigma dW_s \rangle - \int_{\tau_k}^{\sigma_k} \alpha_s^2 \langle \nabla_\theta v(X_s, \theta_s), \nabla_\theta f(X_s, \theta_s) \sigma^{-1} dW_s \rangle. \end{aligned} \quad (3.2)$$

Next step is to treat each term on the right hand side of (3.2) separately. For this purpose, let us first set

$$J_t^{(1)} = \alpha_t \sup_{s \in [0, t]} \|v(X_s, \theta_s)\|.$$

By Theorem A.1 and Proposition 2 of [2] there is some $0 < K < \infty$ (that may change from line to line below) and $0 < q < \infty$ such that for t large enough

$$\begin{aligned}\mathbb{E}|J_t^{(1)}|^2 &\leq K\alpha_t^2 \mathbb{E} \left[1 + \sup_{s \in [0, t]} \|X_s\|^q \right] = K\alpha_t^2 \left[1 + \sqrt{t} \frac{\mathbb{E} \sup_{s \in [0, t]} \|X_s\|^q}{\sqrt{t}} \right] \\ &\leq K\alpha_t^2 [1 + \sqrt{t}] \leq K\alpha_t^2 \sqrt{t}.\end{aligned}$$

Let us now consider $p > 0$ such that $\lim_{t \rightarrow \infty} \alpha_t^2 t^{1/2+2p} = 0$ and for any $\delta \in (0, p)$ define the event $A_{t, \delta} = \left\{ J_t^{(1)} \geq t^{\delta-p} \right\}$. Then we have for t large enough such that $\alpha_t^2 t^{1/2+2p} \leq 1$

$$\mathbb{P}(A_{t, \delta}) \leq \frac{\mathbb{E}|J_t^{(1)}|^2}{t^{2(\delta-p)}} \leq K \frac{\alpha_t^2 t^{1/2+2p}}{t^{2\delta}} \leq K \frac{1}{t^{2\delta}}$$

The latter implies that

$$\sum_{n \in \mathbb{N}} \mathbb{P}(A_{2^n, \delta}) < \infty.$$

Therefore, by Borel-Cantelli lemma we have that for every $\delta \in (0, p)$ there is a finite positive random variable $d(\omega)$ and some $n_0 < \infty$ such that for every $n \geq n_0$ one has

$$J_{2^n}^{(1)} \leq \frac{d(\omega)}{2^{n(p-\delta)}}.$$

Thus for $t \in [2^n, 2^{n+1})$ and $n \geq n_0$ one has for some finite constant $K < \infty$

$$J_t^{(1)} \leq \alpha_{2^n} \sup_{s \in [2^n, 2^{n+1}]} |v(X_s, \theta_s)| \leq K\alpha_{2^{n+1}} \sup_{s \in (0, 2^{n+1}]} |v(X_s, \theta_s)| \leq K \frac{d(\omega)}{2^{(n+1)(p-\delta)}} \leq K \frac{d(\omega)}{t^{p-\delta}}.$$

The latter display then guarantees that for $t \geq 2^{n_0}$ we have with probability one

$$J_t^{(1)} \leq K \frac{d(\omega)}{t^{p-\delta}} \rightarrow 0, \text{ as } t \rightarrow \infty. \quad (3.3)$$

Next we consider the term

$$J_{t,0}^{(2)} = \int_0^t \left\| \alpha'_s v(X_s, \theta_s) + \alpha_s (\mathcal{L}_\theta v(X_s, \theta_s) + \alpha_s \text{tr} [\nabla_\theta f(X_s, \theta_s) \nabla_{x_i} \nabla_\theta v(X_s, \theta_s)]_{i=1}^m) \right\| ds.$$

By the bounds of Theorem A.1 we see that there are constants $0 < K < \infty$ (that may change from line to line) and $0 < q < \infty$ such that

$$\begin{aligned}\mathbb{E}|J_{\infty,0}^{(2)}| &\leq K \int_0^\infty (|\alpha'_s| + \alpha_s^2)(1 + \mathbb{E}\|X_s\|^q) ds \\ &\leq K \int_0^\infty (|\alpha'_s| + \alpha_s^2) ds. \\ &\leq K.\end{aligned}$$

The first inequality follows by Theorem A.1, the second inequality follows by Proposition 1 in [2] and the third inequality follows by Condition 2.1.

The latter display implies that the random variable $J_{\infty,0}^{(2)}$ is finite with probability one, which then implies that there is a finite random variable $\bar{J}_0^{(2)}$ such that

$$J_{t,0}^{(2)} \rightarrow \bar{J}_0^{(2)}, \text{ as } t \rightarrow \infty \text{ with probability one.} \quad (3.4)$$

The last term that we need to consider is the martingale term

$$J_{t,0}^{(3)} = \int_{\tau_k}^{\sigma_k} \alpha_s \langle \nabla_x v(X_s, \theta_s), \sigma dW_s \rangle + \int_{\tau_k}^{\sigma_k} \alpha_s^2 \langle \nabla_\theta v(X_s, \theta_s), \nabla_\theta f(X_s, \theta_s) \sigma^{-1} dW_s \rangle.$$

Notice that the Burkholder-Davis-Gundy inequality and the bounds of Theorem A.1 (doing calculations similar to the ones for the term $J_{t,0}^{(2)}$) give us that for some finite constant $K < \infty$, we have

$$\sup_{t>0} \mathbb{E} \left| J_{t,0}^{(3)} \right|^2 \leq K \int_0^\infty \alpha_s^2 ds < \infty$$

Thus, by Doob's martingale convergence theorem there is a square integrable random variable $\bar{J}^{(3)}$ such that

$$J_{t,0}^{(3)} \rightarrow \bar{J}^{(3)}, \text{ as } t \rightarrow \infty \text{ both almost surely and in } L^2. \quad (3.5)$$

Let us now go back to (3.2). Using the terms $J_t^{(1)}$, $J_{t,0}^{(2)}$ and $J_{t,0}^{(3)}$ we can write

$$\|\Gamma_{k,\eta}\| \leq J_{\sigma_k,\eta}^{(1)} + J_{\tau_k}^{(1)} + J_{\sigma_k,\eta,\tau_k}^{(2)} + \|J_{\sigma_k,\eta,\tau_k}^{(3)}\|$$

The last display together with (3.3), (3.4) and (3.5) imply the statement of the lemma. \square

Lemma 3.3 shows that the function \bar{g} and its first two derivatives are uniformly bounded in θ .

Lemma 3.3. *Assume Conditions 2.1, 2.2 and 2.3. For any $q > 0$, there is a constant K such that*

$$\int_{\mathcal{X}} (1 + |x|^q) \pi(dx) \leq C.$$

In addition we also have that there is a constant $C < \infty$ such that $\sum_{i=0}^2 \|\nabla_\theta^i \bar{g}(\theta)\| \leq C$.

Proof. By Theorem 1 in [3], the density μ of the measure π admits, for any p , a constant C_p such that $\mu(x) \leq \frac{C_p}{1+|x|^p}$. Choosing p large enough that $\int_{\mathcal{X}} \frac{1+|x|^q}{1+|x|^p} dy < \infty$, we then obtain

$$\int_{\mathcal{X}} (1 + |x|^q) \pi(dx) \leq \int_{\mathcal{X}} K_p \frac{1 + |x|^q}{1 + |x|^p} dx \leq C.$$

concluding the proof of the first statement of the lemma. Let us now focus on the second part of the lemma. We only prove the claim for $i = 0$, since due to the bounds in Condition 2.3, the proof for $i = 1, 2$ is the same. By Condition 2.3 and by the first part of the lemma, we have that there exist constants $0 < q, K, C < \infty$ such that

$$\bar{g}(\theta) = \int_{\mathcal{X}} \frac{1}{2} \|f(x, \theta) - f^*(x)\|^2 \pi(dx) \leq K \int_{\mathcal{X}} (1 + |x|^q) \pi(dx) \leq C,$$

concluding the proof of the lemma. \square

Our next goal is to show that if the index k is large enough, then \bar{g} decreases, in the sense of Lemma 3.4.

Lemma 3.4. *Assume Conditions 2.1, 2.2 and 2.3. Suppose that there are an infinite number of intervals $I_k = [\tau_k, \sigma_k]$. There is a fixed constant $\gamma = \gamma(\kappa) > 0$ such that for k large enough, one has*

$$\bar{g}(\theta_{\sigma_k}) - \bar{g}(\theta_{\tau_k}) \leq -\gamma. \quad (3.6)$$

Proof. By Itô's formula we have that

$$\begin{aligned} \bar{g}(\theta_{\sigma_k}) - \bar{g}(\theta_{\tau_k}) &= - \int_{\tau_k}^{\sigma_k} \alpha_s \|\nabla \bar{g}(\theta_s)\|^2 ds + \int_{\tau_k}^{\sigma_k} \alpha_s \langle \nabla \bar{g}(\theta_s), \nabla_\theta f(X_s, \theta_s) \sigma^{-1} dW_s \rangle \\ &+ \int_{\tau_k}^{\sigma_k} \frac{\alpha_s^2}{2} \text{tr} [(\nabla_\theta f(X_s, \theta_s) \sigma^{-1})(\nabla_\theta f(X_s, \theta_s) \sigma^{-1})^T \nabla_\theta \nabla_\theta \bar{g}(\theta_s)] ds \\ &+ \int_{\tau_k}^{\sigma_k} \alpha_s \langle \nabla_\theta \bar{g}(\theta_s), \nabla_\theta \bar{g}(\theta_s) - \nabla_\theta g(X_s, \theta_s) \rangle ds \\ &= \Theta_{1,k} + \Theta_{2,k} + \Theta_{3,k} + \Theta_{4,k}. \end{aligned}$$

Let's first consider $\Theta_{1,k}$. Notice that for all $s \in [\tau_k, \sigma_k]$ one has $\frac{\|\nabla \bar{g}(\theta_{\tau_k})\|}{2} \leq \|\nabla \bar{g}(\theta_s)\| \leq 2\|\nabla \bar{g}(\theta_{\tau_k})\|$. Hence, for sufficiently large k , we have the upper bound:

$$\Theta_{1,k} = - \int_{\tau_k}^{\sigma_k} \alpha_s \|\nabla \bar{g}(\theta_s)\|^2 ds \leq - \frac{\|\nabla \bar{g}(\theta_{\tau_k})\|^2}{4} \int_{\tau_k}^{\sigma_k} \alpha_s ds \leq - \frac{\|\nabla \bar{g}(\theta_{\tau_k})\|^2}{8} \lambda,$$

since Lemma 3.2 proved that $\int_{\tau_k}^{\sigma_k} \alpha_s ds \geq \frac{\lambda}{2}$ for sufficiently large k .

We next address $\Theta_{2,k}$ and show that it becomes small as $k \rightarrow \infty$. First notice that we can trivially write

$$\begin{aligned} \Theta_{2,k} &= \int_{\tau_k}^{\sigma_k} \alpha_s \langle \nabla \bar{g}(\theta_s), \nabla_{\theta} f(X_s, \theta_s) \sigma^{-1} dW_s \rangle = \|\nabla \bar{g}(\theta_{\tau_k})\| \int_{\tau_k}^{\sigma_k} \alpha_s \left\langle \frac{\nabla \bar{g}(\theta_s)}{\|\nabla \bar{g}(\theta_{\tau_k})\|}, \nabla_{\theta} f(X_s, \theta_s) \sigma^{-1} dW_s \right\rangle \\ &= \|\nabla \bar{g}(\theta_{\tau_k})\| \int_{\tau_k}^{\sigma_k} \alpha_s \left\langle \frac{\nabla \bar{g}(\theta_s)}{R_s}, \nabla_{\theta} f(X_s, \theta_s) \sigma^{-1} dW_s \right\rangle. \end{aligned}$$

By Condition 2.3 and Itô isometry we have

$$\begin{aligned} \sup_{t>0} \mathbb{E} \left| \int_0^t \alpha_s \left\langle \frac{\nabla \bar{g}(\theta_s)}{R_s}, \nabla_{\theta} f(X_s, \theta_s) \sigma^{-1} dW_s \right\rangle \right|^2 &\leq 4 \mathbb{E} \int_0^{\infty} \alpha_s^2 \|\nabla_{\theta} f(X_s, \theta_s)\|^2 ds \\ &\leq K \int_0^{\infty} \alpha_s^2 (1 + \mathbb{E} \|X_s\|^q) ds < \infty, \end{aligned}$$

where R_s is defined via (3.1). Hence, by Doob's martingale convergence theorem there is a square integrable random variable M such that $\int_0^t \alpha_s \left\langle \frac{\nabla \bar{g}(\theta_s)}{R_s}, \nabla_{\theta} f(X_s, \theta_s) dW_s \right\rangle \rightarrow M$ both almost surely and in L^2 . The latter statement implies that for a given $\epsilon > 0$ there is k large enough such that almost surely

$$\Theta_{2,k} \leq \|\nabla \bar{g}(\theta_{\tau_k})\| \epsilon.$$

We now consider $\Theta_{3,k}$.

$$\mathbb{E} \left\| \int_0^{\infty} \frac{\alpha_s^2}{2} \text{tr} [(\nabla_{\theta} f(X_s, \theta_s) \sigma^{-1})(\nabla_{\theta} f(X_s, \theta_s) \sigma^{-1})^T \nabla_{\theta} \nabla_{\theta} \bar{g}(\theta_s)] ds \right\| \leq C \int_0^{\infty} \frac{\alpha_s^2}{2} \mathbb{E} (1 + \|X_s\|^q) ds < \infty, \quad (3.7)$$

where we have used Condition 2.3 and Lemma 3.3. Bound (3.7) implies that

$$\int_0^{\infty} \frac{\alpha_s^2}{2} \text{tr} [(\nabla_{\theta} f(X_s, \theta_s) \sigma^{-1})(\nabla_{\theta} f(X_s, \theta_s) \sigma^{-1})^T \nabla_{\theta} \nabla_{\theta} \bar{g}(\theta_s)] ds$$

is finite almost surely, which in turn implies that there is a finite random variable Θ_3^{∞} such that

$$\int_0^t \frac{\alpha_s^2}{2} \text{tr} [(\nabla_{\theta} f(X_s, \theta_s) \sigma^{-1})(\nabla_{\theta} f(X_s, \theta_s) \sigma^{-1})^T \nabla_{\theta} \nabla_{\theta} \bar{g}(\theta_s)] ds \rightarrow \Theta_3^{\infty} \text{ as } t \rightarrow \infty,$$

with probability one. Since Θ_3^{∞} is finite, $\int_{\tau_k}^{\sigma_k} \frac{\alpha_s^2}{2} \text{tr} [(\nabla_{\theta} f(X_s, \theta_s) \sigma^{-1})(\nabla_{\theta} f(X_s, \theta_s) \sigma^{-1})^T \nabla_{\theta} \nabla_{\theta} \bar{g}(\theta_s)] ds \rightarrow 0$ as $k \rightarrow \infty$ with probability one.

Finally, we address $\Theta_{4,k}$. Let us consider the function $G(t, x, \theta) = \alpha_t \langle \nabla_{\theta} \bar{g}(\theta), \nabla_{\theta} g(x, \theta) - \nabla_{\theta} \bar{g}(\theta) \rangle$. The function $G(t, x, \theta)$ satisfies the centering condition (A.1) of Theorem A.1. Therefore, the Poisson equation (A.2) with right hand side $G(t, x, \theta)$ will have a unique smooth solution that grows at most polynomially in x . Let us apply Itô formula to the function $u(t, x, \theta)$ that is solution to this Poisson equation.

$$\begin{aligned} u(\sigma, X_{\sigma}, \theta_{\sigma}) - u(\tau, X_{\tau}, \theta_{\tau}) &= \int_{\tau}^{\sigma} \partial_s u(s, X_s, \theta_s) ds + \int_{\tau}^{\sigma} \mathcal{L}_x u(s, X_s, \theta_s) ds + \int_{\tau}^{\sigma} \mathcal{L}_{\theta} u(s, X_s, \theta_s) ds \\ &\quad + \int_{\tau}^{\sigma} \alpha_s \text{tr} [\nabla_{\theta} f(X_s, \theta_s) \nabla_x \nabla_{\theta} u(s, X_s, \theta_s)] ds \\ &\quad + \int_{\tau}^{\sigma} \langle \nabla_x u(s, X_s, \theta_s), \sigma dW_s \rangle + \int_{\tau}^{\sigma} \alpha_s \langle \nabla_{\theta} u(s, X_s, \theta_s), \nabla_{\theta} f(X_s, \theta_s) \sigma^{-1} dW_s \rangle. \end{aligned}$$

One can write $u(s, x, \theta) = \alpha_s v(x, \theta)$. Using these facts and rearranging the previous Itô formula yields

$$\begin{aligned}\Theta_{4,k} &= \int_{\tau_k}^{\sigma_k} \alpha_s \langle \nabla_{\theta} \bar{g}(\theta_t), \nabla_{\theta} g(X_s, \theta_s) - \nabla_{\theta} \bar{g}(\theta_s) \rangle ds = \int_{\tau_k}^{\sigma_k} \mathcal{L}_x u(s, X_s, \theta_s) ds \\ &= \left[\alpha_{\sigma_k} v(X_{\sigma_k}, \theta_{\sigma_k}) - \alpha_{\tau_k} v(X_{\tau_k}, \theta_{\tau_k}) - \int_{\tau_k}^{\sigma_k} \partial_s \alpha_s v(X_s, \theta_s) ds \right] \\ &\quad - \int_{\tau_k}^{\sigma_k} \alpha_s [\mathcal{L}_{\theta} v(X_s, \theta_s) + \alpha_s \text{tr} [\nabla_{\theta} f(X_s, \theta_s) \nabla_x \nabla_{\theta} v(X_s, \theta_s)]] ds \\ &\quad - \int_{\tau_k}^{\sigma_k} \alpha_s \langle \nabla_x v(X_s, \theta_s), \sigma dW_s \rangle - \int_{\tau_k}^{\sigma_k} \alpha_s \langle \nabla_{\theta} v(X_s, \theta_s), \nabla_{\theta} f(X_s, \theta_s) \sigma^{-1} dW_s \rangle.\end{aligned}$$

Following the exact same steps as in the proof of Lemma 3.2 gives us that $\lim_{k \rightarrow \infty} \|\Theta_{4,k}\| \rightarrow 0$ almost surely.

We now return to $\bar{g}(\theta_{\sigma_k}) - \bar{g}(\theta_{\tau_k})$ and provide an upper bound which is negative. For sufficiently large k , we have that:

$$\begin{aligned}\bar{g}(\theta_{\sigma_k}) - \bar{g}(\theta_{\tau_k}) &\leq -\frac{\|\nabla \bar{g}(\theta_{\tau_k})\|^2}{8} \lambda + \|\Theta_{2,k}\| + \|\Theta_{3,k}\| + \|\Theta_{4,k}\| \\ &\leq -\frac{\|\nabla \bar{g}(\theta_{\tau_k})\|^2}{8} \lambda + \|\nabla \bar{g}(\theta_{\tau_k})\| \epsilon + \epsilon + \epsilon.\end{aligned}$$

Choose $\epsilon = \min\{\frac{\lambda \kappa^2}{32}, \frac{\lambda}{32}\}$. On the one hand, if $\|\nabla \bar{g}(\theta_{\tau_k})\| \geq 1$:

$$\begin{aligned}\bar{g}(\theta_{\sigma_k}) - \bar{g}(\theta_{\tau_k}) &\leq -\frac{\|\nabla \bar{g}(\theta_{\tau_k})\|^2}{8} \lambda + \|\nabla \bar{g}(\theta_{\tau_k})\|^2 \epsilon + \epsilon + \epsilon \\ &\leq -3 \frac{\|\nabla \bar{g}(\theta_{\tau_k})\|^2}{32} \lambda + 2\epsilon \leq -3 \frac{\kappa^2}{32} \lambda + 2 \frac{\kappa^2}{32} \lambda \leq -\frac{\kappa^2}{32} \lambda.\end{aligned}$$

On the other hand, if $\|\nabla \bar{g}(\theta_{\tau_k})\| \leq 1$, then

$$\begin{aligned}\bar{g}(\theta_{\sigma_k}) - \bar{g}(\theta_{\tau_k}) &\leq -\frac{\|\nabla \bar{g}(\theta_{\tau_k})\|^2}{8} \lambda + \epsilon + \epsilon + \epsilon \\ &\leq -\frac{4\kappa^2}{32} \lambda + 3\epsilon \leq -4 \frac{\kappa^2}{32} \lambda + 3 \frac{\kappa^2}{32} \lambda \leq -\frac{\kappa^2}{32} \lambda.\end{aligned}$$

Finally, let $\gamma = \frac{\kappa^2}{32} \lambda$ and the proof of the lemma is complete. □

Lemma 3.5. *Assume Conditions 2.1, 2.2 and 2.3. Suppose that there are an infinite number of intervals $I_k = [\tau_k, \sigma_k]$. There is a fixed constant $\gamma_1 < \gamma$ such that for k large enough,*

$$\bar{g}(\theta_{\tau_k}) - \bar{g}(\theta_{\sigma_{k-1}}) \leq \gamma_1.$$

Proof. First, recall that $\|\nabla \bar{g}(\theta_t)\| \leq \kappa$ for $t \in J_k = [\sigma_{k-1}, \tau_k]$. Similar to before, we have that:

$$\begin{aligned}
\bar{g}(\theta_{\tau_k}) - \bar{g}(\theta_{\sigma_{k-1}}) &= - \int_{\sigma_{k-1}}^{\tau_k} \alpha_s \|\nabla \bar{g}(\theta_s)\|^2 ds + \int_{\sigma_{k-1}}^{\tau_k} \alpha_s \langle \nabla \bar{g}(\theta_s), \nabla_\theta f(X_s, \theta_s) \sigma^{-1} dW_s \rangle \\
&+ \int_{\sigma_{k-1}}^{\tau_k} \frac{\alpha_s^2}{2} \text{tr} [(\nabla_\theta f(X_s, \theta_s) \sigma^{-1})(\nabla_\theta f(X_s, \theta_s) \sigma^{-1})^T \nabla_\theta^2 \bar{g}(\theta_s)] ds \\
&+ \int_{\sigma_{k-1}}^{\tau_k} \alpha_s \langle \nabla_\theta \bar{g}(\theta_s), \nabla_\theta \bar{g}(\theta_s) - \nabla_\theta g(X_s, \theta_s) \rangle ds \\
&\leq \int_{\sigma_{k-1}}^{\tau_k} \alpha_s \langle \nabla \bar{g}(\theta_s), \nabla_\theta f(X_s, \theta_s) \sigma^{-1} dW_s \rangle \\
&+ \int_{\sigma_{k-1}}^{\tau_k} \frac{\alpha_s^2}{2} \text{tr} [(\nabla_\theta f(X_s, \theta_s) \sigma^{-1})(\nabla_\theta f(X_s, \theta_s) \sigma^{-1})^T \nabla_\theta^2 \bar{g}(\theta_s)] ds \\
&+ \int_{\sigma_{k-1}}^{\tau_k} \alpha_s \langle \nabla_\theta \bar{g}(\theta_s), \nabla_\theta \bar{g}(\theta_s) - \nabla_\theta g(X_s, \theta_s) \rangle ds. \tag{3.8}
\end{aligned}$$

The right hand side (RHS) of equation (3.8) converges almost surely to 0 as $k \rightarrow \infty$ as a consequence of similar arguments as given in Lemma 3.4. Indeed, the treatment of the second and third terms on the RHS of (3.8) are exactly the same as in Lemma 3.4. It remains to show that the first term on the RHS of (3.8) converges almost surely to 0 as $k \rightarrow \infty$.

$$\begin{aligned}
\int_{\sigma_{k-1}}^{\tau_k} \alpha_s \langle \nabla \bar{g}(\theta_s), \nabla_\theta f(X_s, \theta_s) \sigma^{-1} dW_s \rangle &= \|\nabla \bar{g}(\theta_{\sigma_{k-1}})\| \int_{\sigma_{k-1}}^{\tau_k} \alpha_s \left\langle \frac{\nabla \bar{g}(\theta_s)}{\|\nabla \bar{g}(\theta_{\sigma_{k-1}})\|}, \nabla_\theta f(X_s, \theta_s) \sigma^{-1} dW_s \right\rangle \\
&= \|\nabla \bar{g}(\theta_{\sigma_{k-1}})\| \int_{\sigma_{k-1}}^{\tau_k} \alpha_s \left\langle \frac{\nabla \bar{g}(\theta_s)}{R_s}, \nabla_\theta f(X_s, \theta_s) \sigma^{-1} dW_s \right\rangle. \tag{3.9}
\end{aligned}$$

As shown in Lemma 3.4, $\int_{\sigma_{k-1}}^{\tau_k} \alpha_s \left\langle \frac{\nabla \bar{g}(\theta_s)}{R_s}, \nabla_\theta f(X_s, \theta_s) \sigma^{-1} dW_s \right\rangle \rightarrow 0$ as $k \rightarrow \infty$ almost surely. Finally, note that $\|\nabla \bar{g}(\theta_{\sigma_{k-1}})\| \leq \kappa$ (except when $\sigma_{k-1} = \tau_k$, in which case the interval J_k is length 0 and hence the integral (3.9) over J_k is 0). Then, $\int_{\sigma_{k-1}}^{\tau_k} \alpha_s \langle \nabla \bar{g}(\theta_s), \nabla_\theta f(X_s, \theta_s) \sigma^{-1} dW_s \rangle \rightarrow 0$ as $k \rightarrow \infty$ almost surely.

Therefore, with probability one, $\bar{g}(\theta_{\tau_k}) - \bar{g}(\theta_{\sigma_{k-1}}) \leq \gamma_1 < \gamma$ for sufficiently large k . \square

Proof of Theorem 2.4. Choose a $\kappa > 0$. First, consider the case where there are a finite number of times τ_k . Then, there is a finite T such that $\|\nabla \bar{g}(\theta_t)\| < \kappa$ for $t \geq T$. Now, consider the other case where there are an infinite number of times τ_k and use Lemmas 3.4 and 3.5. With probability one,

$$\begin{aligned}
\bar{g}(\theta_{\sigma_k}) - \bar{g}(\theta_{\tau_k}) &\leq -\gamma = -\frac{\kappa^2}{32} \lambda, \\
\bar{g}(\theta_{\tau_k}) - \bar{g}(\theta_{\sigma_{k-1}}) &\leq \gamma_1 < \gamma, \tag{3.10}
\end{aligned}$$

for sufficiently large k . Choose a K such that (3.10) holds for $k \geq K$. This leads to:

$$\bar{g}(\theta_{\tau_{n+1}}) - \bar{g}(\theta_{\tau_K}) = \sum_{k=K}^n \left[\bar{g}(\theta_{\sigma_k}) - \bar{g}(\theta_{\tau_k}) + \bar{g}(\theta_{\tau_{k+1}}) - \bar{g}(\theta_{\sigma_k}) \right] \leq \sum_{k=K}^n (-\gamma + \gamma_1) < 0.$$

Let $n \rightarrow \infty$ and then $\bar{g}(\theta_{\tau_{n+1}}) \rightarrow -\infty$. However, we also have that by definition $\bar{g}(\theta) \geq 0$. This is a contradiction, and therefore almost surely there are a finite number of times τ_k .

Consequently, there exists a finite time T (possibly random) such that almost surely $\|\nabla \bar{g}(\theta_t)\| < \kappa$ for $t \geq T$. Since the original $\kappa > 0$ was arbitrarily chosen, this shows that $\|\nabla \bar{g}(\theta_t)\| \rightarrow 0$ as $t \rightarrow \infty$ almost surely. \square

4 Numerical Analysis

We implement the continuous stochastic gradient descent for several application cases and numerically analyze the convergence. Section 4.1 studies continuous stochastic gradient descent for the Ornstein-Uhlenbeck process, which is widely-used in finance, physics, and biology. Section 4.2 studies the multidimensional Ornstein-Uhlenbeck process. Section 4.3 estimates the diffusion coefficient in Burger’s equation with continuous stochastic gradient descent. Burger’s equation is a widely used nonlinear partial differential equation which is important to fluid mechanics, nonlinear acoustics, and gas dynamics. Burger’s equation is extensively used in engineering (especially for aerodynamics). In our final example, Section 4.4, we show how continuous stochastic gradient descent can be used for reinforcement learning.

4.1 Ornstein-Uhlenbeck process

The Ornstein-Uhlenbeck (OU) process $X_t \in \mathbb{R}$ satisfies the stochastic differential equation:

$$dX_t = c(m - X_t)dt + dW_t. \quad (4.1)$$

We use continuous stochastic gradient descent to learn the parameters $\theta = (c, m) \in \mathbb{R}^2$.

For the numerical experiments, we use an Euler scheme with a time step of 10^{-2} . The learning rate is $\alpha_t = \min(\alpha, \alpha/t)$ with $\alpha = 10^{-2}$. We simulate data from (4.1) for a particular θ^* and the stochastic gradient descent attempts to learn a parameter θ_t which fits the data well. θ_t is the statistical estimate for θ^* at time t . If the estimation is accurate, θ_t should of course be close to θ^* . This example can be placed in the form of the original class of equations (1.1) by setting $f(x, \theta) = c(m - x)$ and $f^*(x) = f(x, \theta^*)$.

We study 10,500 cases. For each case, a different θ^* is randomly generated in the range $[1, 2] \times [1, 2]$. For each case, we solve for the parameter θ_t over the time period $[0, T]$ for $T = 10^6$. To summarize:

- For cases $n = 1$ to 10,500
 - Generate a random θ^* in $[1, 2] \times [1, 2]$
 - Simulate a single path of X_t given θ^* and simultaneously solve for the path of θ_t on $[0, T]$

The accuracy of θ_t at times $t = 10^2, 10^3, 10^4, 10^5$, and 10^6 is reported in Table 1. Figures 1 and 2 plots the mean error in percent and mean squared error (MSE) against time. In the table and figures, the “error” is $|\theta_t^n - \theta^{*,n}|$ where n represents the n -th case. The “error in percent” is $100 \times \frac{|\theta_t^n - \theta^{*,n}|}{|\theta^{*,n}|}$. The “mean error in percent” is the average of these errors, i.e. $\frac{100}{N} \sum_{n=1}^N \frac{|\theta_t^n - \theta^{*,n}|}{|\theta^{*,n}|}$.

Error/Time	10^2	10^3	10^4	10^5	10^6
Maximum Error	.604	.2615	.0936	.0349	.0105
99% quantile of error	.368	.140	.0480	.0163	.00542
99.9% quantile of error	.470	.1874	.0670	.0225	.00772
Mean squared error	1.92×10^{-2}	2.28×10^{-3}	2.52×10^{-4}	2.76×10^{-5}	2.90×10^{-6}
Mean Error in percent	7.37	2.497	0.811	0.264	0.085
Maximum error in percent	59.92	20.37	5.367	1.79	0.567
99% quantile of error in percent	25.14	9.07	3.05	1.00	0.323
99.9% quantile of error in percent	34.86	12.38	4.12	1.30	0.432

Table 1: Error at different times for the estimate θ_t of θ^* across 10,500 cases. The “error” is $|\theta_t^n - \theta^{*,n}|$ where n represents the n -th case. The “error in percent” is $100 \times \frac{|\theta_t^n - \theta^{*,n}|}{|\theta^{*,n}|}$.

Finally, we also track the objective function $\bar{g}(\theta_t)$ over time. Figure 3 plots the error $\bar{g}(\theta_t)$ against time. Since the limiting distribution $\pi(x)$ of (4.2) is Gaussian with mean m^* and variance $\frac{1}{2c^*}$, we have that:

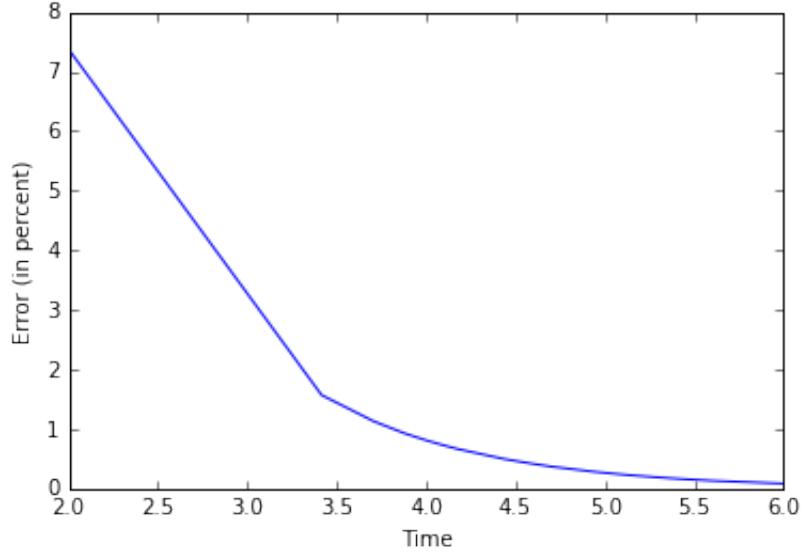


Figure 1: Mean error in percent plotted against time. Time is in log scale.

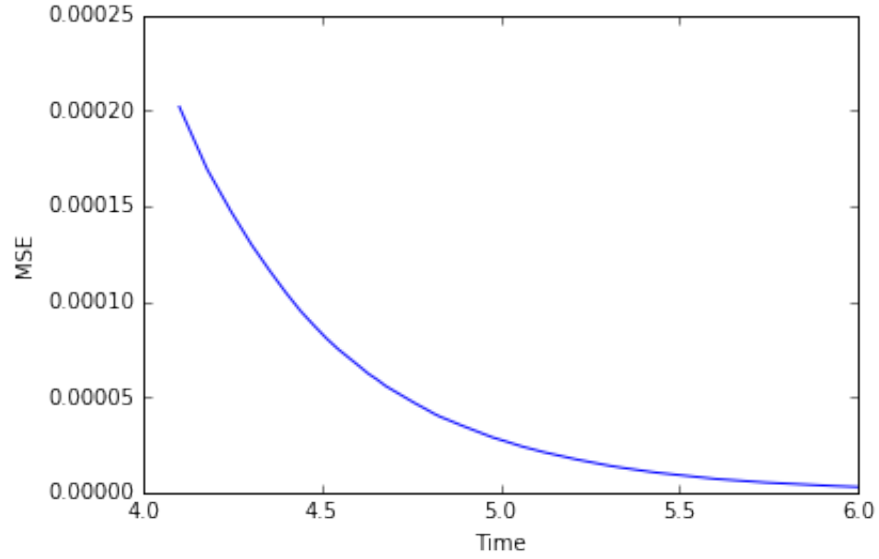


Figure 2: Mean squared error plotted against time. Time is in log scale.

$$\begin{aligned}
\bar{g}(\theta) &= \int \left(c^*(m^* - x) - c(m - x) \right)^2 \pi(x) dx \\
&= (c^*m^* - cm)^2 + (c^* - c)^2 \left(\frac{1}{2c^*} + (m^*)^2 \right) + 2(c^*m^* - cm)(c - c^*)m^*
\end{aligned}$$

4.2 Multidimensional Ornstein-Uhlenbeck process

The multidimensional Ornstein-Uhlenbeck process $X_t \in \mathbb{R}^d$ satisfies the stochastic differential equation:

$$dX_t = (M - AX_t)dt + dW_t. \quad (4.2)$$

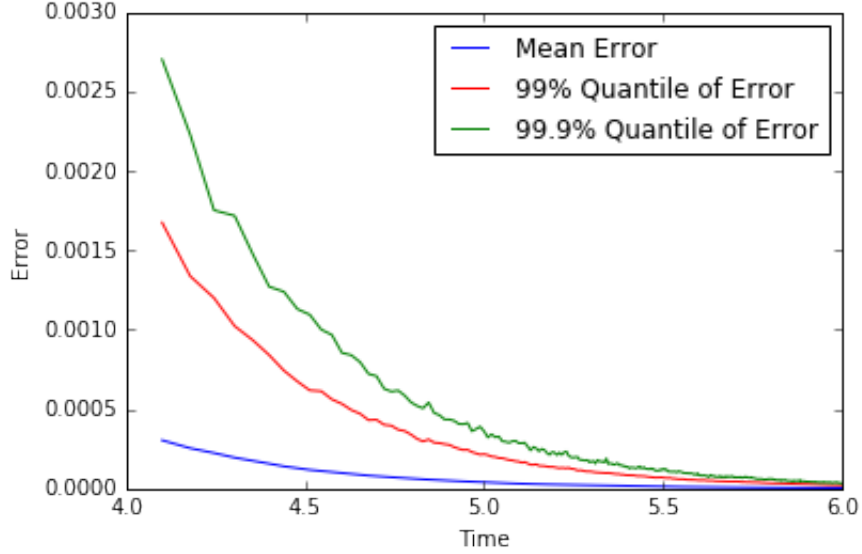


Figure 3: The error $\bar{g}(\theta_t)$ plotted against time. The mean error and the quantiles of the error are calculated from the 10,500 cases. Time is in log scale.

We use continuous stochastic gradient descent to learn the parameters $\theta = (M, A) \in \mathbb{R}^d \times \mathbb{R}^{d \times d}$.

For the numerical experiments, we use an Euler scheme with a time step of 10^{-2} . The learning rate is $\alpha_t = \min(\alpha, \alpha/t)$ with $\alpha = 10^{-1}$. We simulate data from (4.2) for a particular $\theta^* = (M^*, A^*)$ and the stochastic gradient descent attempts to learn a parameter θ_t which fits the data well. θ_t is the statistical estimate for θ^* at time t . If the estimation is accurate, θ_t should of course be close to θ^* . This example can be placed in the form of the original class of equations (1.1) by setting $f(x, \theta) = M - Ax$ and $f^*(x) = f(x, \theta^*)$.

The matrix A^* must be generated carefully to ensure that X_t is ergodic and has a stable equilibrium point. If some of A^* 's eigenvalues have negative real parts, then X_t can become unstable and grow arbitrarily large. Therefore, we randomly generate matrices A^* which are strictly diagonally dominant. A 's eigenvalues are therefore guaranteed to have positive real parts and X_t will be ergodic. To generate random strictly diagonally dominant matrices A^* , we first generate $A_{i,j}^*$ randomly in $[1, 2]$ for $i \neq j$. Then, we set $A_{i,i}^* = \sum_{j \neq i} A_{i,j}^* + U_{i,i}$ where $U_{i,i}$ is generated randomly in $[1, 2]$. M_i^* for $i = 1, \dots, d$ is also generated randomly in $[1, 2]$.

We study 525 cases and analyze the error in Table 2. Figures 4 and 5 plot the error over time.

Error/Time	10^2	10^3	10^4	10^5	10^6
Maximum Error	2.89	.559	.151	.043	.013
99% quantile of error	2.19	.370	.0957	.0294	.00911
99.9% quantile of error	2.57	.481	.118	.0377	.0117
Mean squared error	8.05×10^{-1}	2.09×10^{-2}	1.38×10^{-3}	1.29×10^{-4}	1.25×10^{-5}
Mean Error in percent	34.26	6.18	1.68	0.52	0.161
Maximum error in percent	186.3	41.68	10.98	3.81	1.03
99% quantile of error in percent	109.2	23.9	6.98	2.15	0.657
99.9% quantile of error in percent	141.2	31.24	8.64	2.84	0.879

Table 2: Error at different times for the estimate θ_t of θ^* across 525 cases. The “error” is $|\theta_t^n - \theta^{*,n}|$ where n represents the n -th case. The “error in percent” is $100 \times \frac{|\theta_t^n - \theta^{*,n}|}{|\theta^{*,n}|}$.

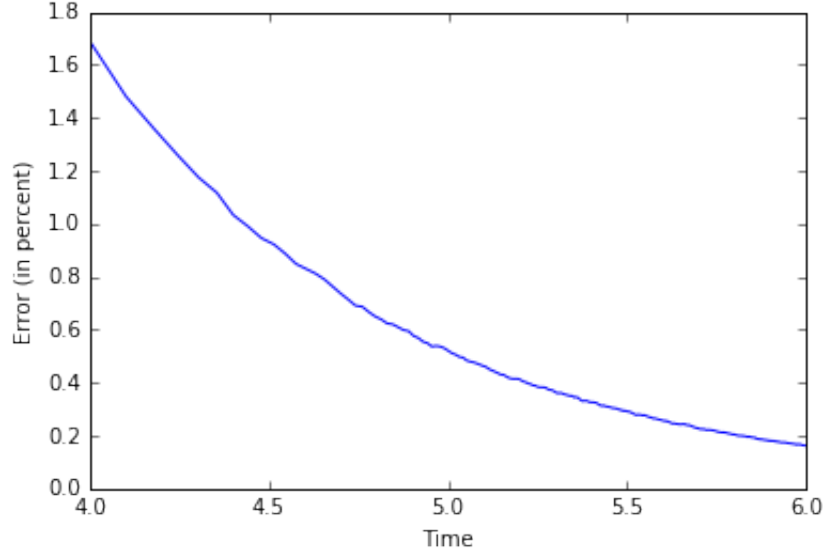


Figure 4: Mean error in percent plotted against time. Time is in log scale.

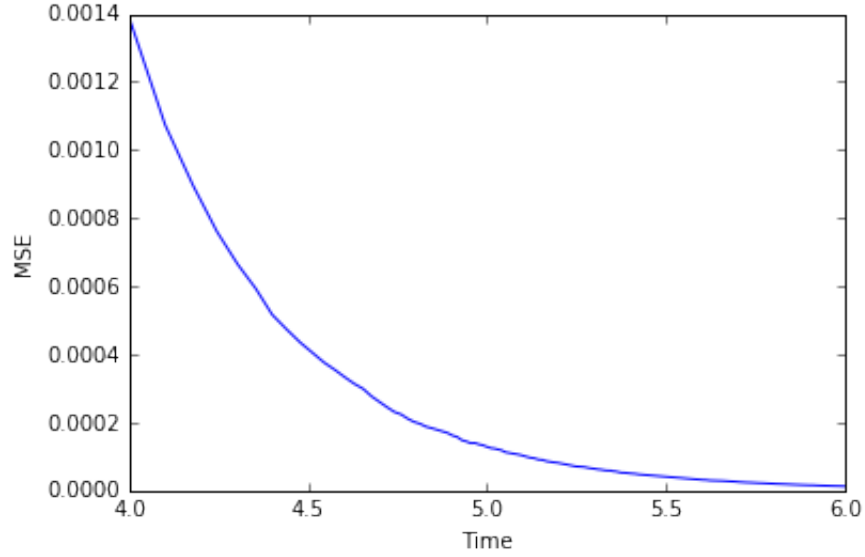


Figure 5: Mean squared error plotted against time. Time is in log scale.

4.3 Burger's Equation

The stochastic Burger's equation that we consider is given by:

$$\frac{\partial u}{\partial t}(t, x) = \theta \frac{\partial^2 u}{\partial x^2} - u(t, x) \frac{\partial u}{\partial x}(t, x) + \sigma \frac{\partial^2 W(t, x)}{\partial t \partial x}, \quad (4.3)$$

where $x \in [0, 1]$ and $W(t, x)$ is a Brownian sheet. The finite-difference discretization of (4.3) satisfies a system of nonlinear stochastic differential equations (for instance, see [4] or [5]). We use continuous stochastic gradient descent to learn the diffusion parameter θ .

We use the following finite difference scheme for Burger's equation:

$$du(t, x_i) = \theta \frac{u(t, x_{i+1}) - 2u(t, x_i) + u(t, x_{i-1}))}{\Delta x^2} dt - u(t, x_i) \frac{u(t, x_{i+1}) - u(t, x_{i-1}))}{2\Delta x} dt + \frac{\sigma}{\sqrt{\Delta x}} dW_t^i, \quad (4.4)$$

For our numerical experiment, the boundary conditions $u(t, x = 0) = 0$ and $u(t, x = 1) = 1$ are used and $\sigma = 0.1$. (4.4) is simulated with the Euler scheme (i.e., we solve Burger’s equation with explicit finite difference). A spatial discretization of $\Delta x = .01$ and a time step of 10^{-5} are used. The learning rate is $\alpha_t = \min(\alpha, \alpha/t)$ with $\alpha = 10^{-3}$. The small time step is needed to avoid instability in the explicit finite difference scheme. Even if $\sigma = 0$, note that an implicit finite difference scheme cannot be used due to the nonlinearity of the partial differential equation (4.3). We simulate data from (4.3) for a particular diffusion coefficient θ^* and the stochastic gradient descent attempts to learn a diffusion parameter θ_t which fits the data well. θ_t is the statistical estimate for θ^* at time t . If the estimation is accurate, θ_t should of course be close to θ^* .

This example can be placed in the form of the original class of equations (1.1). Let f_i be the i -th element of the function f . Then, $f_i(u, \theta) = \theta \frac{u(t, x_{i+1}) - 2u(t, x_i) + u(t, x_{i-1}))}{\Delta x^2} - u(t, x_i) \frac{u(t, x_{i+1}) - u(t, x_{i-1}))}{2\Delta x}$. Similarly, let f_i^* be the i -th element of the function f^* . Then, $f_i^*(u) = f_i(u, \theta^*)$.

We study 525 cases. For each case, a different θ^* is randomly generated in the range $[.1, 10]$. This represents a wide range of physical cases of interest, with θ^* ranging over two orders of magnitude. For each case, we solve for the parameter θ_t over the time period $[0, T]$ for $T = 100$.

The accuracy of θ_t at times $t = 10^{-1}, 10^0, 10^1$, and 10^2 is reported in Table 3. Figures 6 and 7 plots the mean error in percent and mean squared error against time. The convergence of θ_t to θ^* is fairly rapid in time.

Error/Time	10^{-1}	10^0	10^1	10^2
Maximum Error	.1047	.106	.033	.0107
99% quantile of error	.08	.078	.0255	.00835
Mean squared error	1.00×10^{-3}	9.25×10^{-4}	1.02×10^{-4}	1.12×10^{-5}
Mean Error in percent	1.26	1.17	0.4	0.13
Maximum error in percent	37.1	37.5	9.82	4.73
99% quantile of error in percent	12.6	18.0	5.64	1.38

Table 3: Error at different times for the estimate θ_t of θ^* across 525 cases. The “error” is $|\theta_t^n - \theta^{*,n}|$ where n represents the n -th case. The “error in percent” is $100 \times |\theta_t^n - \theta^{*,n}|/|\theta^{*,n}|$.

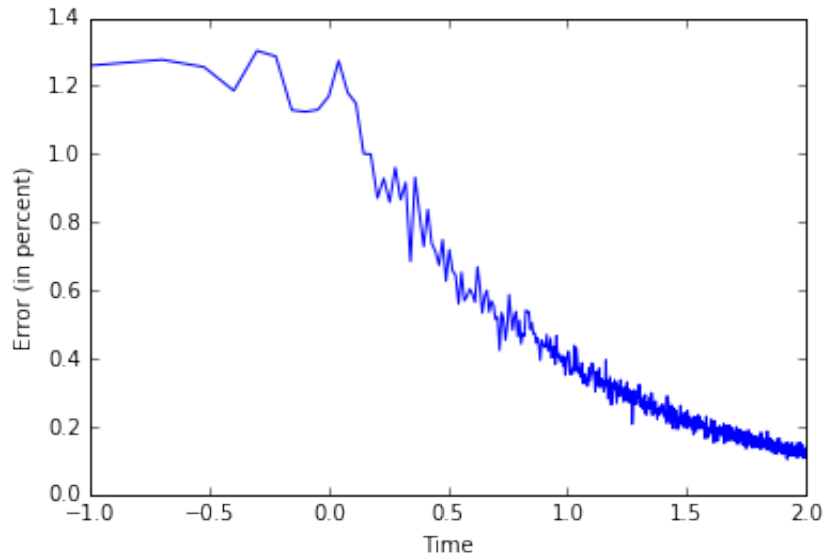


Figure 6: Mean error in percent plotted against time. Time is in log scale.

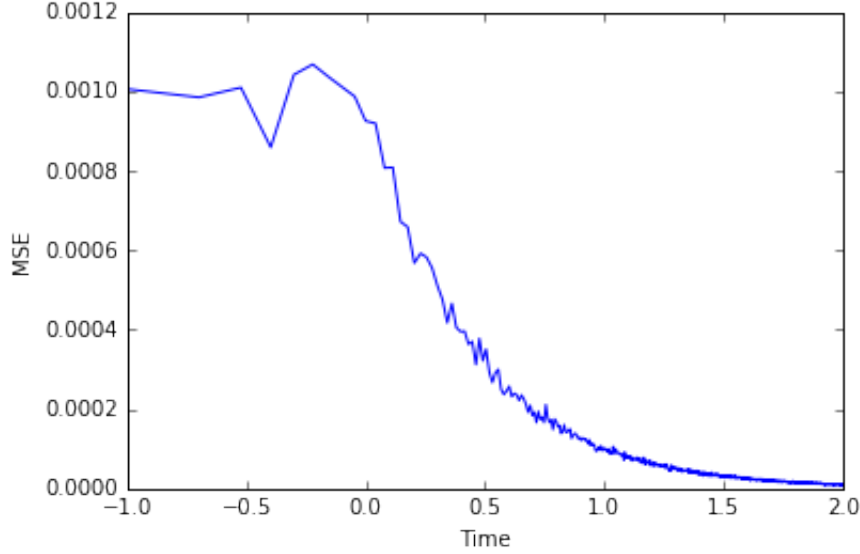


Figure 7: Mean squared error plotted against time. Time is in log scale.

4.4 Reinforcement Learning

We consider the classic reinforcement learning problem of balancing a pole on a moving cart (see [11]). The goal is to balance a pole on a cart and to keep the cart from moving outside the boundaries via applying a force of ± 10 Newtons.

The position x of the cart, the velocity \dot{x} of the cart, angle of the pole β , and angular velocity $\dot{\beta}$ of the pole are observed. The dynamics of $s = (x, \dot{x}, \beta, \dot{\beta})$ satisfy a set of ODEs (see [11]):

$$\begin{aligned}\ddot{\beta}_t &= \frac{g \sin \beta_t + \cos \beta_t \left[\frac{-F_t - ml \dot{\beta}_t^2 \sin \beta_t + \mu_c \operatorname{sgn}(\dot{x}_t)}{m_c + m} \right] - \frac{\mu_p \dot{\beta}_t}{ml}}{l \left[\frac{4}{3} - \frac{m \cos^2 \beta_t}{m_c + m} \right]}, \\ \ddot{x}_t &= \frac{F_t + ml [\dot{\beta}_t^2 \sin \beta_t - \ddot{\beta}_t \cos \beta_t] - \mu_c \operatorname{sgn}(\dot{x}_t)}{m_c + m},\end{aligned}\tag{4.5}$$

where g is the acceleration due to gravity, m_c is the mass of the cart, m is the mass of the pole, $2l$ is the length of the pole, μ_c is the coefficient of friction of the cart on the ground, μ_p is the coefficient of friction of the pole on the cart, and $F_t \in \{-10, 10\}$ is the force applied to the cart.

For this example, $f^*(s) = (\dot{x}, \ddot{x}, \dot{\beta}, \ddot{\beta})$. The model $f(s, \theta) = (f_1(s, \theta), f_2(s, \theta), f_3(s, \theta), f_4(s, \theta))$ where $f_i(s, \theta)$ is a single-layer neural network with rectified linear units.

$$f_i(s, \theta) = W^{2,i} h(W^{1,i} s + b^{1,i}) + b^{2,i},\tag{4.6}$$

where $\theta = \{W^{2,i}, W^{1,i}, b^{1,i}, b^{2,i}\}_{i=1}^4$ and $h(z) = (\sigma(z_1), \dots, \sigma(z_d))$ for $z \in \mathbb{R}^d$. The function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is a rectified linear unit (ReLU): $\sigma(v) = \max(v, 0)$. We learn the parameter θ using stochastic gradient descent.

The boundary is $x = \pm 2.4$ meters and the pole must not be allowed to fall below $\beta = \frac{24}{360\pi}$ radians (the frame of reference is chosen such that the perfectly upright is 0 radians). A reward of +1 is received every 0.02 seconds if $\|x\| \leq 2.4$ and $\|\theta\| \leq \frac{24}{360\pi}$. A reward of -100 is received (and the episode ends) if the cart moves beyond $x = \pm 2.4$ or the pole falls below $\beta = \frac{24}{360\pi}$ radians. The sum of these rewards across the entire episode is the reward for that episode. The initial state $(x, \dot{x}, \beta, \dot{\beta})$ at the start of an episode is generated uniformly at random on $[-.05, .05]^4$. For our numerical experiment, we assume that the rule for receiving the rewards and the distribution of the initial state are both known. An action of ± 10 Newtons may be chosen every 0.02 seconds. This force is then applied for the duration of the next 0.02 seconds.

The goal, of course, is to statistically learn the optimal actions in order to achieve the highest possible reward. This requires both: 1) statistically learning the physical dynamics of $(x, \dot{x}, \beta, \dot{\beta})$ and 2) finding the optimal actions given these dynamics in order to achieve the highest possible reward. The dynamics $(x, \dot{x}, \beta, \dot{\beta})$ satisfy the set of ODEs (4.5); these dynamics can be learned using continuous stochastic gradient descent. We use a neural network for f . Given the estimated dynamics f , we use a policy gradient method to estimate the optimal actions. The approach is summarized below.

- For episodes $0, 1, 2, \dots$:
 - For time $[0, T_{\text{end of episode}}]$:
 - * Update the model $f(s, \theta)$ for the dynamics using continuous stochastic gradient descent.
 - Periodically update the optimal policy $\mu(s, a, \theta^\mu)$ using policy gradient method. The optimal policy is learned using data simulated from the model $f(s, \theta)$. Actions are randomly selected via the policy μ .

The policy μ is a neural network with parameters θ^μ . We use a single hidden layer with rectified linear units followed by a softmax layer for $\mu(s, a, \theta^\mu)$ and train it using policy gradients.² The policy $\mu(s, a, \theta^\mu)$ gives the probability of taking action a conditional on being in the state s .

$$\mathbb{P}[F_t = 10 | s_t = s] = \mu(s_t, 10, \theta^\mu) = \sigma_0(W^2 h(W^1 s + b^1) + b^2), \quad (4.7)$$

where $\sigma_0(v) = \frac{e^v}{1+e^v}$. Of course, $\mathbb{P}[F_t = -10 | s_t = s] = \mu(s_t, -10, \theta^\mu) = 1 - \mu(s_t, 10, \theta^\mu)$.

525 cases are run, each for 25 hours. The optimal policy is learned using the estimated dynamics $f(s, \theta)$ and is updated every 5 episodes. Table 4 reports the results at fixed episodes using continuous stochastic gradient descent. Table 5 reports statistics on the number of episodes required until a target episodic reward (100, 500, 1000) is first achieved.

Reward/Episode	10	20	30	40	45
Maximum Reward	-20	981	2.21×10^4	6.64×10^5	9.22×10^5
90% quantile of reward	-63	184	760	8354	1.5×10^4
Mean reward	-78	67	401	5659	1.22×10^4
10% quantile of reward	-89	-34	36	69	93
Minimum reward	-92	-82	-61	-46	-23

Table 4: Reward at the k -th episode across the 525 cases using continuous stochastic gradient descent to learn the model dynamics.

Alternatively, one could directly apply policy gradient to learn the optimal action using the observed data. This approach does not use continuous stochastic gradient descent to learn the model dynamics, but instead directly learns the optimal policy from the data. Again using 525 cases, we report the results in Table 6 for directly learning the optimal policy without using continuous stochastic gradient descent to learn the model dynamics. Comparing Tables 4 and 6, it is clear that using continuous stochastic gradient descent to learn the model dynamics allows for the optimal policy to be learned significantly more quickly. The rewards are much higher when using continuous stochastic gradient descent (see Table 4) than when not using it (see Table 6).

²Let $r_{e,t}$ be the reward for episode e at time t . Let $R_{t,e} = \sum_{t'=t+1}^{T_{\text{end of episode}}} \gamma^{t'-t} r_{e,t'}$ be the cumulative discounted reward from episode e after time t where $\gamma \in [0, 1]$ is the discount factor. Stochastic gradient descent is used to learn the parameter θ^μ : $\theta^\mu \leftarrow \theta^\mu + \eta_e R_{t,e} \frac{\partial}{\partial \theta^\mu} \log \mu(s_t, a_t, \theta^\mu)$ where η_e is the learning rate. In practice, the cumulative discounted rewards are often normalized across an episode.

Number of episodes/Target reward	100	500	1000
Maximum	39	134	428
90% quantile	23	49	61
Mean	18	34	43
10% quantile	13	21	26
Minimum	11	14	17

Table 5: For each case, we record the number of episodes required until the target reward is first achieved using continuous stochastic gradient descent. Statistics (maximum, quantiles, mean, minimum) for the number of episodes required until the target reward is first achieved.

Reward/Episode	10	20	30	40	100	500	750
Maximum Reward	51	1	15	77	121	1748	1.91×10^5
90% quantile of reward	-52	-48	-42	8354	-11	345	2314
Mean reward	-73	-72	-69	-68	-53	150	1476
10% quantile of reward	-88	-88	-87	69	-83	-1	63
Minimum reward	-92	-92	-92	-92	-92	-81	-74

Table 6: Reward at the k -th episode across the 525 cases using policy gradient to learn the optimal policy.

A On a related Poisson equation

Next we recall the following regularity result from [3] on the Poisson equations in the whole space, appropriately stated to cover our case of interest.

Theorem A.1. *Let Conditions 2.2 and 2.3 be satisfied. Assume that $G(x, \theta) \in C^{\alpha, 2}(\mathcal{X}, \mathbb{R}^n)$,*

$$\int_{\mathcal{X}} G(x, \theta) \pi(dx) = 0, \quad (\text{A.1})$$

and that for some positive constants K and q ,

$$\sum_{i=0}^2 \left| \frac{\partial^i G}{\partial \theta^i}(x, \theta) \right| \leq K (1 + |x|^q)$$

Let \mathcal{L}_x be the infinitesimal generator for the X process. Then the Poisson equation

$$\mathcal{L}_x u(x, \theta) = -G(x, \theta), \quad \int_{\mathcal{X}} u(x, \theta) \pi(dx) = 0 \quad (\text{A.2})$$

has a unique solution that satisfies $u(x, \cdot) \in C^2$ for every $x \in \mathcal{X}$, $\partial_{\theta}^2 u \in C(\mathcal{X} \times \mathbb{R}^n)$ and there exist positive constants K' and q' such that

$$\sum_{i=0}^2 \left| \frac{\partial^i u}{\partial \theta^i}(x, \theta) \right| + \left| \frac{\partial^2 u}{\partial x \partial \theta}(x, \theta) \right| \leq K' (1 + |x|^{q'}).$$

References

- [1] Dimitri P. Bertsekas and John N. Tsitsiklis, Gradient convergence in gradient methods via errors, *SIAM Journal of Optimization*, Vol.10, No. 3, (2000), pp. 627-642.
- [2] E. Pardoux and A.Yu. Veretennikov, On Poisson equation and diffusion approximation I, *Annals of Probability*, Vol. 29, No. 3, (2001), pp. 1061-1085.

- [3] E. Pardoux and A. Y. Veretennikov, On Poisson equation and diffusion approximation 2, *The Annals of Probability* 31 (3) (2003) 1166–1192.
- [4] A. Davie and J. Gaines, Convergence of numerical schemes for the solution of parabolic stochastic partial differential equations, *Mathematics of Computation*, Vol. 70, No. 233, (2001), pp. 121-134.
- [5] A. Alabert and I. Gyongy, On numerical approximation of stochastic Burger’s equation, *From stochastic calculus to mathematical finance*, Springer Berling Heidelberg, 2006. 1–15.
- [6] H. Kushner and G. Yin, *Stochastic Approximation and Recursive Algorithms and Applications*, Second Edition. *Springer*, 2003.
- [7] A. Benveniste, M. Metivier, and P. Priouret, *Adaptive Algorithms and Stochastic Approximations*. *Springer-Verlag*, 2012.
- [8] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning. Book in preparation for MIT Press*, 2016.
- [9] A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro, Robust stochastic approximation to stochastic programming, *SIAM Journal of Optimization*, Vol.19, No. 4, (2009), pp. 1574-1609.
- [10] M. Raginsky and J. Boudrie, Continuous-time stochastic mirror descent on a network: variance reduction, consensus, convergence, *IEEE Conference on Decision and Control*, 2012.
- [11] A. Barto, R. Sutton, and C. Anderson, Neuronlike Adaptive Elements that can solve difficult learning control problem, *IEEE Transactions on Systems, Man, and Cybernetics*, Vol.5, (1983), pp. 834-846.