

Implicit stochastic approximation

Panos Toulis^α and Edoardo M. Airolidi^α

^αDepartment of Statistics, Harvard University,
Cambridge, MA, 02138, USA

October 14, 2015

Abstract

The need to carry out parameter estimation from massive data has reinvigorated interest in iterative estimation methods, in statistics and machine learning. Classic work includes deterministic gradient-based methods, such as quasi-newton, and stochastic gradient descent and its variants, including adaptive learning rates, acceleration and averaging. Current work increasingly relies on methods that employ proximal operators, leading to updates defined through implicit equations, which need to be solved at each iteration. Such methods are especially attractive in modern problems with massive data because they are numerically stable and converge with minimal assumptions, among other reasons. However, while the majority of existing methods can be subsumed into the gradient-free stochastic approximation framework developed by [Robbins and Monro \(1951\)](#), there is no such framework for methods with implicit updates. Here, we conceptualize a gradient-free implicit stochastic approximation procedure, and develop asymptotic and non-asymptotic theory for it. This new framework provides a theoretical foundation for gradient-based procedures that rely on implicit updates, and opens the door to iterative estimation methods that do not require a gradient, nor a fully known likelihood.

Contents

1	Introduction	3
2	Implicit stochastic approximation	4
3	Theory of implicit stochastic approximation	5
3.1	Convergence of implicit stochastic approximation	6
3.2	Non-asymptotic analysis	6
3.3	Asymptotic distribution	8
4	Algorithms for implicit stochastic approximation	8
5	Application in parameter estimation	9
5.1	Likelihood-based estimation	10
5.2	Estimation with likelihood known up to normalizing constant	11
5.3	Likelihood-free estimation	12
6	Conclusion	14
A	Appendix	19
A.1	Implicit stochastic approximation	19
A.2	Theory of implicit stochastic approximation	19
A.2.1	Convergence	20
A.3	Non-asymptotic analysis	21
A.4	Asymptotic distribution	27

1 Introduction

Robbins and Monro (1951) considered the problem of estimating the zero θ_* of a function $h : \mathbb{R}^p \rightarrow \mathbb{R}$. Specifically, for every fixed $\theta \in \mathbb{R}^p$, they assumed that the exact value $h(\theta)$ is unknown but can be unbiasedly estimated by a random variable W_θ , such that $\mathbb{E}(W_\theta) = h(\theta)$. Starting from an estimate θ_0^{rm} , Robbins and Monro (1951) recursively approximated θ_* as follows,

$$\theta_n^{\text{rm}} = \theta_{n-1}^{\text{rm}} - \gamma_n W_{\theta_{n-1}^{\text{rm}}}, \quad (1)$$

where $\{\gamma_n\}$ is a decreasing sequence of positive numbers, known as the learning rate sequence; typically, $\gamma_n \propto 1/n$ such that $\sum \gamma_i^2 < \infty$ to guarantee convergence, and $\sum \gamma_i = \infty$ to guarantee that convergence can be towards any point in \mathbb{R}^p . Robbins and Monro (1951) proved convergence in quadratic mean of procedure (1) —also known as the “Robbins-Monro procedure”—, i.e., $\mathbb{E}(\|\theta_n^{\text{rm}} - \theta_*\|^2) \rightarrow 0$. Since then, several other authors have explored its theoretical properties, for example, Ljung et al. (1992); Kushner and Yin (2003); Borkar (2008). Due to its remarkable simplicity and computational efficiency, the Robbins-Monro procedure has found widespread applications across scientific disciplines, including statistics (Nevel’son et al., 1973), engineering (Benveniste et al., 1990), and optimization (Nesterov, 2004).

Recently, the Robbins-Monro procedure has received renewed interest for its applicability in parameter estimation with large data sets, and its connections to stochastic optimization methods. In such settings, even though there is a finite data set D , the Robbins-Monro procedure can be applied with W_θ being, for example, the log-likelihood of θ at a single data point that is sampled with replacement from D . In this case, the theory of Robbins and Monro (1951) implies that $\mathbb{E}_D(\|\theta_n^{\text{rm}} - \hat{\theta}_n\|^2) \rightarrow 0$, where the expectation is now with respect to the empirical distribution of data points in D , and $\hat{\theta}_n$ is an estimator of θ_* , such as the maximum-likelihood estimator, or maximum a-posteriori if regularization is used. In this work, we will study the theoretical properties of a modified stochastic approximation method in the more general setting with a stream of data points, but will also discuss applications to iterative estimation with a finite data set.

Despite its remarkable properties, a well-known issue with the Robbins-Monro procedure is that the sequence $\{\gamma_n\}$ crucially affects both its numerical stability and convergence. Regarding stability, consider a simple example where θ_n^{rm} in Eq.(1) is approximating a scale parameter, and thus needs to be positive. Clearly, even if the previous iterate θ_{n-1}^{rm} is positive, there is no guarantee that the update in Eq.(1) will provide a positive next iterate θ_n^{rm} . Regarding convergence, the Robbins-Monro procedure can be arbitrarily slow if γ_n is even slightly misspecified. For example, let $\gamma_n = \gamma_1/n$, and assume there exists the *scalar potential* $H : \mathbb{R}^p \rightarrow \mathbb{R}$, such that $\nabla H(\theta) = h(\theta)$, for all θ . If H is strongly convex with parameter c , then $\mathbb{E}(\|\theta_n^{\text{rm}} - \theta_*\|^2) = O(n^{-\epsilon})$ if $\epsilon = 2c\gamma_1 < 1$ (Nemirovski et al., 2009, Section 1); (Moulines and Bach, 2011, Section 3.1). In summary, large learning rates can make the iterates θ_n^{rm} diverge numerically, whereas small rates can make the iterates converge very slowly to θ_* . Importantly, these conflicting requirements for stability and fast convergence are very hard to reconcile in practice, especially in high-dimensional problems (Toulis and Airoldi, 2015a, Section 3.5), which renders the

Robbins-Monro procedure, and all its derived methods, inapplicable without heuristic modifications.

2 Implicit stochastic approximation

Our idea to improve the Robbins-Monro procedure (1) is to transform its iteration into a stochastic fixed-point equation as follows,

$$\theta_n^{\text{im}} = \theta_{n-1}^{\text{im}} - \gamma_n W_{\theta_n^{(*)}}, \text{ s.t.} \quad (2)$$

$$\theta_n^{(*)} = \theta_{n-1}^{\text{im}} - \gamma_n h(\theta_n^{(*)}) \stackrel{\text{def}}{=} \mathbb{E}(\theta_n^{\text{im}} | \mathcal{F}_{n-1}), \quad (3)$$

where \mathcal{F}_{n-1} is the σ -algebra adapted to $\{\theta_0^{\text{im}}, \theta_1^{\text{im}}, \dots, \theta_{n-1}^{\text{im}}\}$. The *implicit stochastic approximation* method of Eq.(2) and Eq.(3) can be motivated as the limit of a sequence of improved Robbins-Monro procedures, as follows. First, for convenience, drop the superscript *rm* from procedure (1), and fix the sample history \mathcal{F}_{n-1} . Then, $\mathbb{E}(\theta_n | \mathcal{F}_{n-1}) = \theta_{n-1} - \gamma_n h(\theta_{n-1}) \triangleq \theta_n^{(1)}$. Assuming that $\theta_n^{(1)}$ can be computed even though h is unknown, then it is reasonable to expect that it is better to use $\theta_n^{(1)}$ instead of θ_{n-1} when drawing W_θ in procedure (1). This leads to update $\theta_n = \theta_{n-1} - \gamma_n W_{\theta_n^{(1)}}$. In turn, this update has expected value $\theta_{n-1} - \gamma_n h(\theta_n^{(1)}) \triangleq \theta_n^{(2)}$, and thus $\theta_n^{(2)}$ can be used instead of $\theta_n^{(1)}$, and so on. This argument can be repeated ad infinitum producing the following sequence of improved Robbins-Monro procedures,

$$\begin{aligned} \theta_n &= \theta_{n-1} - \gamma_n W_{\theta_{n-1}}, \\ \theta_n &= \theta_{n-1} - \gamma_n W_{\theta_n^{(1)}}, \\ \theta_n &= \theta_{n-1} - \gamma_n W_{\theta_n^{(2)}}, \\ &\dots \\ \theta_n &= \theta_{n-1} - \gamma_n W_{\theta_n^{(*)}}, \end{aligned}$$

where $\theta_n^{(i)}$ is defined recursively as $\theta_n^{(i)} = \theta_{n-1} - \gamma_n h(\theta_n^{(i-1)})$. The initial iterate, $\theta_n^{(0)} = \theta_{n-1}$, corresponds to the classic Robbins-Monro procedure (1). The limit iterate, $\theta_n^{(\infty)} = \theta_n^{(*)}$, is the fixed point of Eq.(3), and corresponds to implicit stochastic approximation (2).

The improvement achieved by the fixed-point equation (3) can be explained through the following simple argument. First, take norms in Eq.(3) to obtain

$$\|\theta_{n-1}^{\text{im}} - \theta_\star\|^2 \geq \|\theta_n^{(*)} - \theta_\star\|^2 + 2\gamma_n h(\theta_n^{(*)})^\top (\theta_n^{(*)} - \theta_\star).$$

As before, suppose that the scalar potential H exists and is convex. Then, it follows $h(\theta)^\top (\theta - \theta_\star) \geq 0$, which implies that $\|\theta_n^{(*)} - \theta_\star\|^2 \leq \|\theta_{n-1}^{\text{im}} - \theta_\star\|^2$. Let $\xi_n = W_{\theta_n^{(*)}} - h(\theta_n^{(*)})$, such that $\mathbb{E}(\xi_n | \mathcal{F}_{n-1}) = 0$, then Eq.(2) can be written as $\theta_n^{\text{im}} = \theta_n^{(*)} - \gamma_n \xi_n$. Thus, the fixed-point equation contracts $\theta_n^{(*)}$ towards θ_\star , and potentially contracts θ_n^{im} as well. This idea is also closely related to *proximal operators* in optimization because Eq.(3) can be written as

$$\theta_n^{(*)} = \arg \min_{\theta} \left\{ \frac{1}{2} \|\theta - \theta_{n-1}^{\text{im}}\|^2 + \gamma_n H(\theta) \right\}. \quad (4)$$

The right-hand side of Eq.(4) is a *proximal operator* (Bauschke and Combettes, 2011), mapping θ_{n-1}^{im} to the intermediate point $\theta_n^{(*)}$, and has θ_* as a fixed point. Interest on proximal operators has surged in recent years because they are non-expansive and converge with minimal assumptions. Furthermore, they can be applied on non-smooth objectives, and can easily be combined in modular algorithms for optimization in large-scale and distributed settings (Parikh and Boyd, 2013).

On the practical side, the key advantage of the contractive property of implicit stochastic approximation (2) compared to classic stochastic approximation (1) is that it is no longer required to have small learning rates γ_n for stability. This resolves the conflicting requirements for stability and convergence in classic stochastic approximation, and provides valuable flexibility in choosing the learning rate γ_n . The theoretical results of the following section do confirm that implicit stochastic approximation converges under minimal conditions, maintains the asymptotic properties of classic stochastic approximation, but also has remarkable stability in short-term iterations.

On the downside, the implementation of implicit stochastic approximation is generally a challenging task because h is unknown. In specific, the point $\theta_n^{(*)}$ cannot be computed from Eq.(3) because function h is unknown, otherwise we could simply set $\theta_n^{\text{im}} = \theta_n^{(*)}$ in Eq.(2). However, as mentioned before, Eq.(2) can be written as $\theta_n^{\text{im}} = \theta_n^{(*)} - \gamma_n \xi_n$, where $\mathbb{E}(\xi_n | \mathcal{F}_{n-1}) = 0$, and therefore we only need a noisy estimate of the intermediate point $\theta_n^{(*)}$ in order to implement implicit stochastic approximation. We leverage this result to provide several concrete implementations in Section 4.

3 Theory of implicit stochastic approximation

The symbol $\|\cdot\|$ denotes the L_1 vector/matrix norm. We also define the error random variables $\xi_n \triangleq W_{\theta_{n-1}^{(*)}} - h(\theta_{n-1}^{(*)})$, such that $\mathbb{E}(\xi_n | \mathcal{F}_{n-1}) = 0$. The parameter space for θ will be \mathbb{R}^p without loss of generality. For a positive scalar sequence a_n , the sequence $b_n = O(a_n)$ is such that $b_n \leq ca_n$, for some fixed $c > 0$, and every n ; the sequence $b_n = o(a_n)$ is such that $b_n/a_n \rightarrow 0$ in the limit where $n \rightarrow \infty$. $b_n \downarrow 0$ denotes a positive sequence decreasing towards zero. We further assume that implicit stochastic approximation (2) operates under a combination of the following assumptions.

Assumption 1. *It holds, $\gamma_n = \gamma_1 n^{-\gamma}$, $\gamma_1 > 0$ and $\gamma \in (1/2, 1]$.*

Assumption 2. *The regression function h is Lipschitz with parameter L , i.e., for all θ_1, θ_2 ,*

$$\|h(\theta_1) - h(\theta_2)\| \leq L \|\theta_1 - \theta_2\|.$$

Assumption 3. *Function h satisfies either*

- (a) $(\theta - \theta_*)^\top h(\theta) > 0$, for all θ , or
- (b) $(\theta_n^{(*)} - \theta_*)^\top h(\theta_n^{(*)}) \geq \delta_n \|\theta_n^{(*)} - \theta_*\|^2$, where $\delta_n = \delta_1 n^{-\delta}$, $\delta_1 > 0$ and $0 < \gamma + \delta \leq 1$, for all n .

Assumption 4. *There exists a scalar potential $H : \mathbb{R}^p \rightarrow \mathbb{R}$ such that $\nabla H(\theta) = h(\theta)$, for all θ .*

Assumption 5. *There exists fixed $\sigma^2 > 0$ such that, for every n ,*

$$\mathbb{E}(\|\xi_n\|^2 | \mathcal{F}_{n-1}) \leq \sigma^2.$$

Assumption 6. *Let $\Xi_n \stackrel{\text{def}}{=} \mathbb{E}(\xi_n \xi_n^\top | \mathcal{F}_{n-1})$, then $\|\Xi_n - \Xi\| = O(1)$, and $\|\Xi_n - \Xi\| \rightarrow 0$ for fixed positive-definite matrix Ξ . Furthermore, if $\sigma_{n,s}^2 = \mathbb{E}(\mathbb{I}_{\|\xi_n\|^2 \geq s/\gamma_n} \|\xi_n\|^2)$, then for all $s > 0$, $\sum_{i=1}^n \sigma_{i,s}^2 = o(n)$ if $\gamma_n \propto n^{-1}$, or $\sigma_{n,s}^2 = o(1)$ otherwise.*

Assumption 3(a) is a typical convexity assumption. Assumption 3(b) is stronger than the convexity assumption, but weaker than the assumption of “strong convexity” that is typical in the literature. Assumption 5 was first introduced by Robbins and Monro (1951), and has since been standard in stochastic approximation analysis. Assumption 6 is the Lindeberg condition that is used to prove asymptotic normality of θ_n^{im} . Overall, our assumptions are weaker than the assumptions used in classic stochastic approximation; compare, for example, Assumptions 1-6 with assumptions (A1)-(A4) of Borkar (2008, Section 2.1), or assumptions of Benveniste et al. (1990, Theorem 15).

3.1 Convergence of implicit stochastic approximation

In Theorem 1 we derive a proof of almost-sure convergence of implicit stochastic approximation, which relies on the supermartingale lemma of Robbins and Siegmund (1985).

Theorem 1. *Suppose that Assumptions 1, 2, 3(a), and 5 hold. Then the iterates θ_n^{im} of the implicit stochastic approximation procedure (2) converge almost-surely to θ_* ; i.e., $\theta_n^{\text{im}} \rightarrow \theta_*$, such that $h(\theta_*) = 0$, almost-surely.*

The conditions for almost-sure convergence of implicit stochastic approximation are weaker than the conditions for classic stochastic approximation. For example, to show almost-sure convergence for standard stochastic approximation methods, it is typically assumed that the iterates θ_n^{rm} are almost-surely bounded (Borkar, 2008, Assumption (A4)).

3.2 Non-asymptotic analysis

In this section, we prove results on upper bounds for the deviance $\mathbb{E}(H(\theta_n^{\text{im}}) - H(\theta_*))$ and the errors $\mathbb{E}(\|\theta_n^{\text{im}} - \theta_*\|^2)$. This provides information on the rate of convergence, as well as the stability of implicit stochastic approximation methods. Theorem 2 on deviance uses Assumption 3(a), which only assumes non-strong convexity of H , whereas Theorem 3 on squared errors uses Assumption 3(b), which is weaker than strong convexity.

Theorem 2. Suppose that Assumptions 1, 2, 3(a), 4, and 5 hold. Define $\Gamma^2 \triangleq \mathbb{E}(\|\theta_0^{\text{im}} - \theta_\star\|^2) + \sigma^2 \sum_{i=1}^\infty \gamma_i^2 + \gamma_1^2 \sigma^2$. Then, if $\gamma \in (2/3, 1]$, there exists $n_{0,1} < \infty$ such that, for all $n > n_{0,1}$,

$$\mathbb{E}(H(\theta_n^{\text{im}}) - H(\theta_\star)) \leq \left[\frac{2\Gamma^2}{\gamma\gamma_1} + o(1) \right] n^{-1+\gamma}.$$

If $\gamma \in (1/2, 2/3)$, there exists $n_{0,2} < \infty$ such that, for all $n > n_{0,2}$,

$$\mathbb{E}(H(\theta_n^{\text{im}}) - H(\theta_\star)) \leq \left[\Gamma\sigma\sqrt{L\gamma_1} + o(1) \right] n^{-\gamma/2}.$$

Otherwise, $\gamma = 2/3$ and there exists $n_{0,3} < \infty$ such that, for all $n > n_{0,3}$,

$$\mathbb{E}(H(\theta_n^{\text{im}}) - H(\theta_\star)) \leq \left[\frac{3 + \sqrt{9 + 4\gamma_1^3 L\sigma^2/\Gamma^2}}{2\gamma_1/\Gamma^2} + o(1) \right] n^{-1/3}.$$

There are two main results in Theorem 2. First, the rates of convergence for the deviance are either $O(n^{-1+\gamma})$ or $O(n^{-\gamma/2})$, depending on the learning rate parameter γ . These rates match standard stochastic approximation results under non-strong convexity; see, for example, Theorem 4 of [Moulines and Bach \(2011\)](#). Second, there is a uniform decay of the expected deviance towards zero, since the constants $n_{0,1}, n_{0,2}, n_{0,3}$ can be made small, depending on the desired accuracy in the constants of the upper-bounds in Theorem 2. In contrast, in standard stochastic approximation methods under non-strong convexity, there is a term $\exp(4L^2\gamma_1^2 n^{1-2\gamma})$ ([Moulines and Bach, 2011](#), Theorem 4), which can amplify the initial conditions arbitrarily. Thus, implicit stochastic approximation has similar asymptotic properties to classic stochastic approximation, but is significantly more stable.

Theorem 3. Suppose that Assumptions 1, 3(b), and 5 hold, and define $\zeta_n \triangleq \mathbb{E}(\|\theta_n^{\text{im}} - \theta_\star\|^2)$ and $\kappa \triangleq 1 + 2\gamma_1\delta_1$. Then, if $\gamma + \delta < 1$, for every $n > 0$ it holds,

$$\zeta_n \leq e^{-\log \kappa \cdot n^{1-\gamma-\delta}} \zeta_0 + \sigma^2 \frac{\gamma_1 \kappa}{\delta_1} n^{-\gamma+\delta} + O(n^{-\gamma+\delta-1}).$$

Otherwise, if $\gamma = 1, \delta = 0$, it holds,

$$\zeta_n \leq e^{-\log \kappa \cdot \log n} \zeta_0 + \sigma^2 \frac{\gamma_1 \kappa}{\delta_1} n^{-1} + O(n^{-2}).$$

There are two main results in Theorem 3. First, if H is strongly convex ($\delta = 0$), then the rate of convergence of $\mathbb{E}(\|\theta_n^{\text{im}} - \theta_\star\|^2)$ is $O(n^{-\gamma})$, which matches the rate of convergence for classic stochastic approximation under strong convexity ([Benveniste et al., 1990](#), Theorem 22, p.244). Second, there is an exponential discounting of initial conditions ζ_0 regardless of the specification of the learning rate parameter γ_1 and the Lipschitz parameter L . In stark contrast, in classic stochastic approximation there exists a term $\exp(L^2\gamma_1^2 n^{1-2\gamma})$ in front of the initial conditions ζ_0 , which can make the approximation diverge numerically if γ_1 is misspecified with respect to the Lipschitz parameter L ([Moulines and Bach, 2011](#), Theorem 1). Thus, as in the non-strongly convex case of Theorem 2, implicit stochastic approximation has similar asymptotic rates to classic stochastic approximation, but is also significantly more stable.

3.3 Asymptotic distribution

Asymptotic distributions are well-studied in classic stochastic approximation. In this section, we leverage this theory to show that iterates from implicit stochastic approximation are asymptotically normal. The following theorem establishes this result using Theorem 1 of [Fabian \(1968\)](#); see also ([Ljung et al., 1992](#), Chapter II.8).

Theorem 4. *Suppose that Assumptions 1, 2, 3(a), 5, and 6 hold. Suppose also that $(2\gamma_1 J_h(\theta_*) - I)$ is positive-definite, where $J_h(\theta)$ is the Jacobian of h at θ , and I is the $p \times p$ identity matrix. Then, the iterate θ_n^{im} of implicit stochastic approximation (2) is asymptotically normal, such that*

$$n^{\gamma/2}(\theta_n^{\text{im}} - \theta_*) \rightarrow \mathcal{N}_p(0, \Sigma).$$

The covariance matrix Σ is the unique solution of

$$(\gamma_1 J_h(\theta_*) - I/2)\Sigma + \Sigma(\gamma_1 J_h(\theta_*)^\top - I/2) = \Xi.$$

The asymptotic distribution of the implicit iterate θ_n^{im} is identical to the asymptotic distribution of θ_n^{rm} , as derived by [Fabian \(1968\)](#). Intuitively, in the limit, $\theta_n^{(*)} \approx \theta_{n-1}^{\text{im}} + O(\gamma_n)$ with high probability, and thus implicit stochastic approximation behaves like the classic one. We also note that if Ξ commutes with $J_h(\theta_*)$ then Σ can be derived in closed-form as $\Sigma = (2\gamma_1 J_h(\theta_*) - I)^{-1}\Xi$.

4 Algorithms for implicit stochastic approximation

In the previous section, we showed that implicit stochastic approximation has similar asymptotic properties to classic stochastic approximation, but it is also more stable. Therefore, implicit stochastic approximation is arguably a superior form of stochastic approximation.

However, the main drawback of implicit stochastic approximation is that the intermediate value $\theta_n^{(*)}$, which is necessary to compute the update in Eq.(3), is not readily available because it depends on the regression function $h(\cdot)$ that is unknown. Thus, every implementation of implicit stochastic approximation needs to approximate or estimate $\theta_n^{(*)}$ at every iteration.

One strategy for such implementation is to approximate $\theta_n^{(*)}$ through a separate standard stochastic approximation procedure; i.e., at every n th iteration of procedure (2) run a Robbins-Monro procedure x_k for $k = 1, 2, \dots$, starting from $x_0 = \theta_{n-1}^{\text{im}}$ and updating as follows,

$$x_k = x_{k-1} - a_k(\gamma_n W_{x_{k-1}} + x_{k-1} - \theta_{n-1}^{\text{im}}), \quad (5)$$

where $a_k \propto 1/k$ satisfies the Robbins-Monro conditions, and $W_{x_{k-1}}$ are independent draws with fixed $\theta = x_{k-1}$. By the properties of the Robbins-Monro procedure, x_k converges to x_∞ that satisfies $\gamma_n h(x_\infty) + x_\infty - \theta_{n-1}^{\text{im}} = 0$, and thus $x_k \rightarrow \theta_n^{(*)}$, by Eq.(3). For practical

reasons, it would be sufficient to apply algorithm (5) for a finite number of steps, say K , even before x_k has converged to $\theta_n^{(*)}$, and then use x_K in place of $\theta_n^{(*)}$ in procedure (2). We note that Eq.(5) is in fact a stochastic fixed-point iteration, where one critical necessary condition is that the mapping on the right-hand side of Eq.(5) is non-expansive (Borkar, 2008, Theorem 4, Section 10.3); this condition might be hard to validate in practice.

Another implementation is possible if the analytic form of W_θ is known, even though the regression function $h(\theta)$ is not. The idea is, rather counter-intuitively, to use the next iterate θ_n^{im} as an *estimator* of $\theta_n^{(*)}$. This is reasonable because, by definition, $\mathbb{E}(\theta_n^{\text{im}}|\mathcal{F}_{n-1}) = \theta_n^{(*)}$, i.e., θ_n^{im} is an unbiased estimator of $\theta_n^{(*)}$. Thus, procedure (2) could be approximately implemented as follows,

$$\theta_n^{\text{im}} = \theta_{n-1}^{\text{im}} - \gamma_n W_{\theta_n^{\text{im}}}. \quad (6)$$

The next iterate θ_n^{im} appears on both sides of Eq.(6) and thus Eq.(6) is a multi-dimensional fixed-point equation, also known as *implicit update*. Algorithm (6) appears to be the most practical implementation of implicit stochastic approximation, which justifies why the implicit update of algorithm (6) lends its name to implicit stochastic approximation.

Now, the main problem in applying algorithm (6) is to solve efficiently the fixed-point equation at every iteration and calculate θ_n^{im} . To solve this problem we distinguish two cases. First, let $W_\theta \equiv s(\theta)U$, where $s(\theta)$ is a scalar random variable that may depend on parameter θ , and U is unit-norm vector random variable that does not depend on θ . Then, it is easy to see that algorithm (6) is equivalent to

$$\theta_n^{\text{im}} = \theta_{n-1}^{\text{im}} - \gamma_n \lambda_n s(\theta_{n-1}^{\text{im}})U, \quad (7)$$

where the scalar λ_n satisfies

$$\lambda_n s(\theta_{n-1}^{\text{im}}) = s(\theta_{n-1}^{\text{im}} - \gamma_n \lambda_n s(\theta_{n-1}^{\text{im}})U). \quad (8)$$

Therefore, at the n th iteration we can first draw U , then solve the one-dimensional fixed-point equation (8), given that the analytic form of $s(\theta)$ is known, and finally calculate θ_n^{im} through Eq.(7). The solution of the one-dimensional fixed-point equation is computationally efficient in a large family of statistical models when W_θ is the negated gradient of the log-likelihood; this family includes, for example, generalized linear models and convex M-estimation (Toulis and Airoldi, 2015a, Theorem 4.1). Second, consider the more general case $W_\theta = s(\theta)U_\theta$, where U_θ is a random vector that can depend on parameter θ . The problem is now is that we cannot sample $U_{\theta_n^{\text{im}}}$ because it depends on the next iterate. However, we can simply draw $U_{\theta_{n-1}^{\text{im}}}$ instead, which reduces to the previous case where $W_\theta = s(\theta)U$, and therefore algorithm (7) is again applicable.

5 Application in parameter estimation

A key application of stochastic approximation is in iterative estimation of the parameters of a statistical model (Nevel'son et al., 1973, Chapter 8); (Ljung et al., 1992, Chapter 10);

(Borkar, 2008, Section 10.2). In particular, consider a stream of i.i.d. data points (X_n, Y_n) , $n = 1, 2, \dots$, where the outcome $Y_n \in \mathbb{R}^d$ is distributed conditional on covariates $X_n \in \mathbb{R}^p$ according to known density $f(Y_n; X_n, \theta_*)$, but unknown model parameters $\theta_* \in \mathbb{R}^p$. We will consider two cases, one where the likelihood is fully known and can be calculated easily at each data point (X_n, Y_n) , and one where the likelihood is not fully known (e.g., it is known up to normalizing constant) and there is a finite data set instead of a stream. We will also discuss a third case of likelihood-free estimation.

5.1 Likelihood-based estimation

Consider the random variable $W_\theta \stackrel{\text{def}}{=} -\nabla \log f(Y_n; X_n, \theta)$, where the regression function $h(\theta) \stackrel{\text{def}}{=} \mathbb{E}(W_\theta)$ is still unknown. In many statistical models the log-likelihood $\log f$ is concave (Bickel and Doksum, 2015), and thus θ_* is the unique point for which $h(\theta_*) \stackrel{\text{def}}{=} -\mathbb{E}(\nabla \log f(Y_n; X_n, \theta_*)) = 0$. Therefore the standard stochastic approximation procedure (1) can be written as

$$\theta_n^{\text{rm}} = \theta_{n-1}^{\text{rm}} + \gamma_n \nabla \log f(Y_n; X_n, \theta_{n-1}^{\text{rm}}). \quad (9)$$

Stochastic approximation theory implies $\theta_n^{\text{rm}} \rightarrow \theta_*$ and therefore θ_n^{rm} is a consistent estimator of θ_* . Procedure (9) is known as *stochastic gradient descent* (SGD) in optimization and signal processing (Coraluppi and Young, 1969), and has been fundamental in modern machine learning with large data sets (Zhang, 2004; Bottou, 2010; Toulis and Airoldi, 2015b).

An implicit stochastic approximation version of procedure (9) can be implemented through algorithm (6), resulting in the iteration

$$\theta_n^{\text{im}} = \theta_{n-1}^{\text{im}} + \gamma_n \nabla \log f(Y_n; X_n, \theta_n^{\text{im}}). \quad (10)$$

Procedure (10) is known as the *incremental proximal method* in optimization (Bertsekas, 2011), or as *implicit stochastic gradient descent* in statistical estimation (Toulis et al., 2014), and has shown superior performance to standard stochastic gradient descent, both in theory and applications (Toulis et al., 2014; Toulis and Airoldi, 2015a). In particular, in accordance to the theoretical properties of their stochastic approximation counterparts, implicit SGD has identical asymptotic efficiency and convergence rate as standard SGD, but it is significantly more stable. The stochastic proximal gradient algorithm (Singer and Duchi, 2009; Rosasco et al., 2014) is related to implicit SGD. However, in contrast to implicit SGD, the proximal gradient algorithm first makes a standard update (forward step), and then an implicit update (backward step), which may increase convergence speed, but may also introduce instability due to the forward step.

Example. Let $\theta_* \in \mathbb{R}$ be the true parameter of a normal model with i.i.d. observations $Y_n|X_n \sim \mathcal{N}(X_n\theta_*, 1)$, $X_n, Y_n \in \mathbb{R}$. Thus, $\log f(Y_n; X_n, \theta) = -\frac{1}{2}(Y_n - X_n\theta)^2$, and $\nabla \log f(Y_n; X_n, \theta) = (Y_n - X_n\theta)X_n$. Assume $\gamma_n = \gamma_1/n$ as the learning rate. Then, the SGD procedure (9) is given by

$$\theta_n^{\text{rm}} = (1 - \gamma_n X_n^2) \theta_{n-1}^{\text{rm}} + \gamma_n Y_n X_n. \quad (11)$$

Procedure (11) is known as the least mean squares filter (LMS) in signal processing, or as the Widrow-Hoff algorithm (Widrow and Hoff, 1960). The implicit SGD procedure for this problem, using update (10), is given by

$$\theta_n^{\text{im}} = \frac{1}{1 + \gamma_n X_n^2} \theta_{n-1}^{\text{im}} + \frac{\gamma_n}{1 + \gamma_n X_n^2} Y_n X_n. \quad (12)$$

Procedure (12) is also known as the normalized least mean squares filter (NLMS) in signal processing (Nagumo and Noda, 1967). From Eq. (11) we see that it is crucial for standard SGD to have a well-specified learning rate parameter γ_1 . For instance, assume fixed $X_n^2 = x^2$ for simplicity, then if $\gamma_1 x^2 \gg 1$ the iterate θ_n^{rm} will diverge to a value $O(2^{\gamma_1 x^2} / \sqrt{\gamma_1 x^2})$. In contrast, a very large γ_1 will not cause divergence in implicit SGD, but it will simply put more weight on the n th observation $Y_n X_n$ than the previous iterate θ_{n-1}^{im} . Moreover, from a statistical perspective, implicit SGD specifies the correct averaging by weighing the estimate and observation according to the inverse of information $(1 + \gamma_n X_n^2)$.

5.2 Estimation with likelihood known up to normalizing constant

We now consider the case where the likelihood is known up to a normalizing constant, which arises frequently in practice (Gelman and Meng, 1998). In such situations, a large family of estimation methods relies on being able to sample from the underlying model, e.g., through Metropolis-Hastings, which does not require knowledge of the normalization constant. For example, the method of Markov chain Monte Carlo maximum-likelihood estimation (Geyer, 1991) uses simulations to estimate ratios of normalization constants, which appear in the maximization of the log-likelihood over a finite data set.

The same idea underlying simulation-based methods can be applied in iterative estimation with stochastic approximation. Consider a finite data set with N data points, and a complete sufficient statistic S_n that can be calculated on the n th data point. Let $\hat{S}(\theta; k)$ denote the averaged value of the statistic over k independent data points that are simulated conditional on θ . Then, the iteration

$$\begin{aligned} n &\sim \text{Uniform}\{1, 2, \dots, N\} \\ \theta_n &= \theta_{n-1} + \gamma_n \left(S_n - \hat{S}(\theta_{n-1}; k) \right), \end{aligned} \quad (13)$$

can converge to θ_* under typical conditions of stochastic approximation. The learning rates γ_n can be chosen as $\gamma_n = \gamma_1/n$, as it is common, and adaptive schemes are possible (Lai and Robbins, 1979). The constant k affects the variance of the stochastic part in update (13): higher k lowers the variance, but also increases the computational burden because more simulations are needed. Resolution of such statistical/computational trade-offs typically require problem-specific considerations.

The implicit counterpart of procedure (13) uses an update as follows,

$$\theta_n = \theta_{n-1} + \gamma_n \left(S_n - \lambda_n \hat{S}(\theta_{n-1}; k) \right), \quad (14)$$

where the scalar λ_n satisfies

$$\lambda_n \hat{S}(\theta_{n-1}; k) = \hat{S}(\theta_n; k) \stackrel{\text{def}}{=} \hat{S}(\theta_{n-1} + \gamma_n(S_n - \lambda_n \hat{S}(\theta_{n-1}; k)); k). \quad (15)$$

Solution of Eq.(15) is particularly challenging because one needs to find the scalar λ_n such that the update (14) leads to a new iterate θ_n with simulated statistics satisfying (15). One idea is to consider only $\lambda_n \in [0, 1]$, which is true under strong convexity; for example, in the normal linear model of Eq.(11), $\lambda_n = 1/(1 + \gamma_n X_n^2)$. More generally, this idea is reasonable because λ_n usually acts as a shrinkage factor (Toulis and Airoldi, 2015a, Section 5). Then, we can repeat simulations on a grid of $m + 1$ values $[0, 1]_m \triangleq \{0, 1/m, 2/m, \dots, 1\}$, and set λ_n as follows,

$$\lambda_n = \arg \min_{\lambda \in [0, 1]_m} \|\lambda \hat{S}(\theta_{n-1}; k) - \hat{S}(\theta_n(\lambda); k)\|^2, \quad (16)$$

where $\theta_n(\lambda) \triangleq \theta_{n-1} + \gamma_n(S_n - \lambda \hat{S}(\theta_{n-1}; k))$. Procedure (14) uses $m \cdot k$ simulations per iteration from the underlying model, and is thus more computationally demanding than procedure (13), which uses only k simulations per iteration. Improvements in computation are possible if certain simulated moments are reused, similar to the idea of Bartz et al. (2009). This is reasonable because the grid of values $[0, 1]_m$ will increasingly yield similar simulated moments as the iteration counter n increases and $\gamma_n \rightarrow 0$.

Example. Suppose we observe N graphs G_n , $n = 1, 2, \dots, N$, and we want to fit an exponential random graph model (ERGM) (Frank and Strauss, 1986; Pattison and Wasserman, 1999; Robins et al., 2007) with density of a graph G defined as

$$f(G; \theta) = \exp(\theta^\top X_G) / c(\theta), \quad (17)$$

where $X_G \in \mathbb{R}^p$ are properties of G , such as number of triangles, number of edges, etc., and $c(\theta)$ is an appropriate normalizing constant. Typically, $c(\theta)$ is hard to compute, and thus remains unknown.

In this case, it is still possible to use the stochastic approximation procedure (18), where we define $S_n \triangleq X_{G_n}$, and $\hat{S}(\theta_{n-1}; k) \triangleq (1/k) \sum_{i=1}^k X_{\tilde{G}_i}$, where \tilde{G}_i is an independent sample from the ERGM model (17) with parameter value fixed at θ_{n-1} . Obtaining such samples is computationally straightforward (Snijders, 2002; Bartz et al., 2009). The implicit procedure (14) implemented through (16) has potential stability advantages, but a more efficient implementation remains an interesting open problem ahead.

5.3 Likelihood-free estimation

An important class of estimation methods in statistics does not rely on likelihood, such as the method of moments, or non-parametric methods. Typically, in such cases, there is a statistic S_n that can be calculated at every n th data point, and the regression function $T(\theta) \triangleq \mathbb{E}(S_n | \theta_\star = \theta)$ is known. Then, the procedure

$$\theta_n = \theta_{n-1} + \gamma_n (S_n - T(\theta_{n-1})), \quad (18)$$

can converge to θ_* under typical stochastic approximation conditions (e.g., convexity of vector-valued function T). If the system of equations $\mathbb{E}(S_n) = T(\theta_*)$ is over-specified, convexification or regularization could be applied, similar, for example, to the way the generalized method of moments (Hall, 2005) resolves over-specification in moment conditions. Interestingly, the form of procedure (18), where the update depends on a discrepancy $S_n - T(\theta_{n-1})$ between an observed statistic and its expected value, also appears in likelihood-based estimation with SGD, e.g., in exponential family models —cf. Eq.(11). The implicit counterpart of procedure (18) has straightforward implementations through the algorithms of Section 4, which we illustrate in the following example.

Example. In their seminal paper, Robbins and Monro (1951) described an application of procedure (1) in iterative quantile estimation. In particular, consider a random variable Z with a cumulative distribution function F . An experimenter wants to find the point θ_* such that $F(\theta_*) = \alpha$, for a fixed $\alpha \in (0, 1)$. The experimenter cannot draw samples of Z , but has access to the random variable $W_\theta = \mathbb{I}\{Z \leq \theta\} - \alpha$, for any value of θ . Robbins and Monro (1951) argued that the procedure

$$\theta_n = \theta_{n-1} - \gamma_n W_{\theta_{n-1}}, \quad (19)$$

will converge to θ_* . Indeed, $\mathbb{E}(W_{\theta_*}) = \mathbb{E}(\mathbb{I}\{Z \leq \theta_*\}) - \alpha = F(\theta_*) - \alpha \stackrel{\text{def}}{=} 0$. Convexity or concavity of F is not required because F is nondecreasing and $F'(\theta_*) > 0$, which are the original conditions for convergence established by Robbins and Monro (1951).

Quantile estimation through implicit stochastic approximation can be accomplished through algorithm (5), with the n th iteration defined as

$$\begin{aligned} x_k &= x_{k-1} - a_k(\gamma_n W_{x_{k-1}} + x_{k-1} - \theta_{n-1}^{\text{im}}), k = 1, 2, \dots, K, \\ \theta_n^{\text{im}} &= x_K, \end{aligned} \quad (20)$$

where $x_0 = \theta_{n-1}^{\text{im}}$, $\alpha_k \propto 1/k$, and K is a small constant; we give concrete values to these constants in a numerical simulation at the end of this example.

The benefits of the implicit procedure over the classic one can be emphatically illustrated through the following numerical exercise. Suppose Z is a standard normal random variable, and $\alpha = 0.999$. In this case, $F(\theta_*) = 0.999$, and thus $\theta_* \approx 3.09$ is in a region of the parameter space where F is very flat; in fact, the region $[2.87, +\infty]$ is within ± 0.1 from the value 0.999 of $F(\theta_*)$. It is well-known that the convergence of the classic Robbins-Monro procedure depends on $F'(\theta_*)$: if the product $\gamma_1 F'(\theta_*)$ is small then convergence will be slow (Toulis and Airolidi, 2015b, Section 2.1). In this example, because $F'(\theta_*) \approx 0.0034$ is very small in the flat region, the learning rate parameter γ_1 needs to compensate by being large, otherwise convergence to θ_* will be very slow. In fact, standard theory suggests that it is optimal to select $\gamma_1 = 1/F'(\theta_*) \approx 294$. However, if γ_1 is large the classic procedure (19) will overshoot early on to a point $\theta_*^+ \gg \theta_*$ such that $F(\theta_*^+) \approx 1$. Return from that point will be very slow because the values of $W_{\theta_{n-1}}$ will be close to zero, since $W_{\theta_{n-1}} \approx 1 - \alpha = 0.001$. For example, suppose $\gamma_1 = 294$ and $\theta_*^+ = 10$,

then to return back to $\theta_\star \approx 3.09$ we will need m steps for which

$$\begin{aligned} \sum_{n=n_0+1}^{n_0+m} \gamma_n W_{\theta_{n-1}} &\approx (10 - 3.09) \\ \gamma_1 \sum_{n=n_0+1}^{n_0+m} 0.001/n &\approx 6.91 \\ \log m &\geq 6.91 \cdot 10^3 / \gamma_1 \approx 23. \end{aligned} \tag{21}$$

The classic procedure will therefore require at least e^{23} steps to return back after this simple overshoot from $\theta_\star = 3.09$ to $\theta_\star^+ = 10$. Smaller values of the learning rate, e.g., $\gamma_1 = 20$, will exacerbate this problem.

In stark contrast, the implicit procedure (20) can neither overshoot nor undershoot. The update (20) satisfies approximately $\theta_n^{\text{im}} + \gamma_n(F(\theta_n^{\text{im}}) - \alpha) \approx \theta_{n-1}^{\text{im}}$, which can be written as $\theta_n^{\text{im}} + \gamma_n F(\theta_n^{\text{im}}) \approx \theta_{n-1}^{\text{im}} + \gamma_n F(\theta_\star)$. Because F is nondecreasing it follows that if $\theta_{n-1}^{\text{im}} < \theta_\star$, then $\theta_n^{\text{im}} > \theta_{n-1}^{\text{im}}$. Similarly, if $\theta_{n-1}^{\text{im}} > \theta_\star$ then $\theta_n^{\text{im}} < \theta_{n-1}^{\text{im}}$. Furthermore, a large value of γ_n will in fact push θ_n^{im} more towards θ_\star . The implicit equation therefore bestows a remarkable stability property to the underlying stochastic approximation procedure.

To validate our theory, we ran 20 replications of both the Robbins-Monro procedure (19) and the implicit procedure (20) for quantile estimation. Both procedures had a learning rate $\gamma_n = \gamma_1/n$, where we picked multiple values of γ_1 to test the stability of the procedures. For the implicit procedure we also set $K = 50$ and $a_k = 10/k$. Each procedure was initialized at $\theta_0 = -10$ and was run for 5,000 iterations, at each replication. The final parameter estimate θ_n , $n = 5,000$, was stored for each procedure and replication. The results are shown in Figure 1.

We observe that, as the learning rate constant γ_1 increases, the classic RM procedure produces iterates that overshoot and get stuck. This is consistent with the theoretical analysis of this section, where we showed that the classic procedure cannot converge after overshooting because the true parameter value θ_\star is in a region where the objective function F is flat. In stark contrast, the implicit procedure remains significantly robust, with iterates around the true parameter value across different learning rates. Importantly, the implicit procedure was only approximately implemented with algorithm (5), where the inner stochastic approximation for x_k was executed for only $K = 40$ iterations within each n th iteration of the main implicit procedure (20).

6 Conclusion

The need to estimate parameters of a statistical model from massive amounts of data has reinvigorated interest in approximate inference procedures. While most gradient-based procedures of current interest can be seen as special case of a gradient-free stochastic approximation method developed by [Robbins and Monro \(1951\)](#), there is no general gradient-free method that can accommodate procedures with updates defined through implicit equations. Here, we conceptualize a gradient-free implicit stochastic approximation

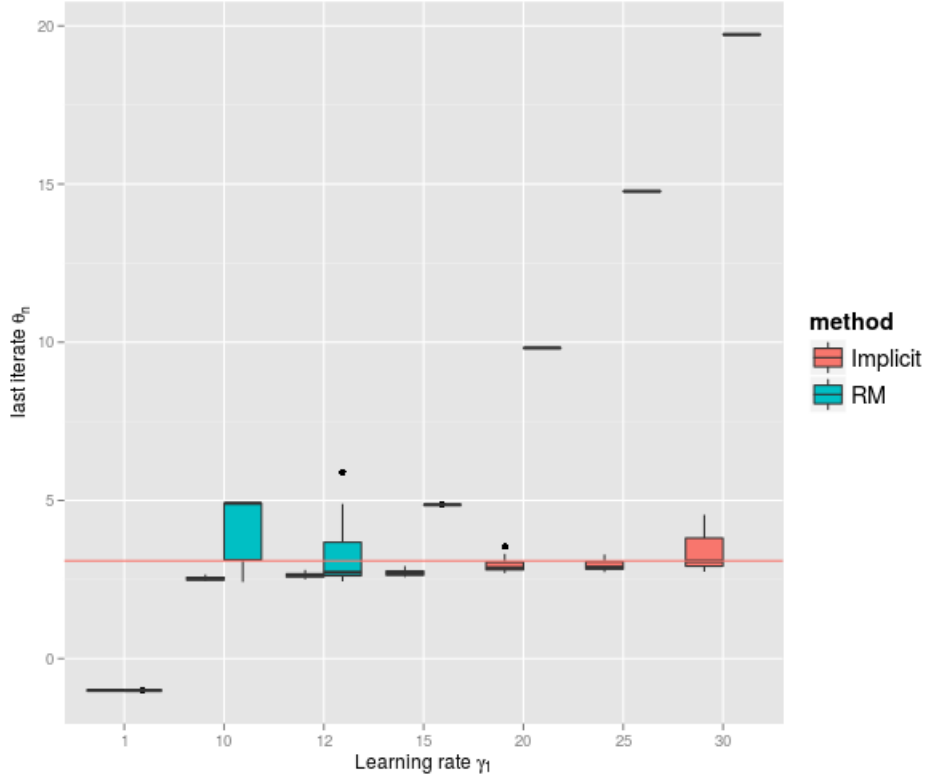


Figure 1: Boxplots of 20 replications of the Robbins-Monro (RM) (19) and the implicit procedure (20). Each replication yields one estimate θ_n , where $n = 5,000$. The true parameter value θ_* is depicted as a red horizontal line at $y = \Phi(0.999) \approx 3.09$, and both procedures start from $\theta_0 = -10$. When the learning rate constant is small (e.g., $\gamma_1 = 0.1$) both procedures are very slow and their iterates are far from the true value. As γ_1 increases, the classic RM procedure produces iterates that overshoot and get stuck (flat boxplots). In contrast, the implicit procedure remains robust, with final iterates around θ_* .

method, and develop asymptotic and non-asymptotic theory for it. This new approximation method provides the theoretical basis for gradient-based procedures that rely on proximal operators (implicit updates), and opens the door to new iterative estimation procedures that do not require access to a gradient or a fully-known likelihood function.

References

- Bartz, K., J. Liu, and J. Blitzstein (2009). Monte carlo maximum likelihood for exponential random graph models: From snowballs to umbrella densities.
- Bauschke, H. H. and P. L. Combettes (2011). *Convex analysis and monotone operator theory in Hilbert spaces*. Springer Science & Business Media.

- Benveniste, A., P. Priouret, and M. Métivier (1990). Adaptive algorithms and stochastic approximations.
- Bertsekas, D. P. (2011). Incremental proximal methods for large scale convex optimization. *Mathematical programming* 129(2), 163–195.
- Bickel, P. J. and K. A. Doksum (2015). *Mathematical Statistics: Basic Ideas and Selected Topics, volume I*, Volume 117. CRC Press.
- Borkar, V. S. (2008). Stochastic approximation. *Cambridge Books*.
- Bottou, L. (2010). Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT'2010*, pp. 177–186. Springer.
- Coraluppi, G. and T. Y. Young (1969). Stochastic signal representation. *Circuit Theory, IEEE Transactions on* 16(2), 155–161.
- Fabian, V. (1968). On asymptotic normality in stochastic approximation. *The Annals of Mathematical Statistics*, 1327–1332.
- Frank, O. and D. Strauss (1986). Markov graphs. *Journal of the american Statistical association* 81(395), 832–842.
- Gelman, A. and X.-L. Meng (1998). Simulating normalizing constants: From importance sampling to bridge sampling to path sampling. *Statistical science*, 163–185.
- Geyer, C. J. (1991). Markov chain monte carlo maximum likelihood.
- Gladyshev, E. (1965). On stochastic approximation. *Theory of Probability & Its Applications* 10(2), 275–278.
- Hall, A. R. (2005). *Generalized method of moments*. Oxford University Press Oxford.
- Kushner, H. J. and G. Yin (2003). *Stochastic approximation and recursive algorithms and applications*, Volume 35. Springer Science & Business Media.
- Lai, T. L. and H. Robbins (1979). Adaptive design and stochastic approximation. *The annals of Statistics*, 1196–1221.
- Ljung, L., G. Pflug, and H. Walk (1992). Stochastic approximation and optimization of random systems.
- Moulines, E. and F. R. Bach (2011). Non-asymptotic analysis of stochastic approximation algorithms for machine learning. In *Advances in Neural Information Processing Systems*, pp. 451–459.
- Nagumo, J.-I. and A. Noda (1967). A learning method for system identification. *Automatic Control, IEEE Transactions on* 12(3), 282–287.

- Nemirovski, A., A. Juditsky, G. Lan, and A. Shapiro (2009). Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization* 19(4), 1574–1609.
- Nesterov, Y. (2004). *Introductory lectures on convex optimization*, Volume 87. Springer Science & Business Media.
- Nevel’son, M. B., R. Z. Khas’minskii, and B. Silver (1973). *Stochastic approximation and recursive estimation*. American Mathematical Society Providence, RI.
- Parikh, N. and S. Boyd (2013). Proximal algorithms. *Foundations and Trends in optimization* 1(3), 123–231.
- Pattison, P. and S. Wasserman (1999). Logit models and logistic regressions for social networks: Ii. multivariate relations. *British Journal of Mathematical and Statistical Psychology* 52(2), 169–194.
- Robbins, H. and S. Monro (1951). A stochastic approximation method. *The annals of mathematical statistics*, 400–407.
- Robbins, H. and D. Siegmund (1985). A convergence theorem for non negative almost supermartingales and some applications. In *Herbert Robbins Selected Papers*, pp. 111–135. Springer.
- Robins, G., P. Pattison, Y. Kalish, and D. Lusher (2007). An introduction to exponential random graph (p^*) models for social networks. *Social networks* 29(2), 173–191.
- Rosasco, L., S. Villa, and B. C. Vũ (2014). Convergence of stochastic proximal gradient algorithm. *arXiv preprint arXiv:1403.5074*.
- Singer, Y. and J. C. Duchi (2009). Efficient learning using forward-backward splitting. In *Advances in Neural Information Processing Systems*, pp. 495–503.
- Snijders, T. A. (2002). Markov chain monte carlo estimation of exponential random graph models. *Journal of Social Structure* 3(2), 1–40.
- Toulis, P., E. Airoldi, and J. Rennie (2014). Statistical analysis of stochastic gradient methods for generalized linear models. In *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, pp. 667–675.
- Toulis, P. and E. M. Airoldi (2015a). Implicit stochastic gradient descent for principled estimation with large datasets. *arXiv preprint arXiv:1408.2923*.
- Toulis, P. and E. M. Airoldi (2015b). Scalable estimation strategies based on stochastic approximations: classical results and new insights. *Statistics and computing* 25(4), 781–795.
- Widrow, B. and M. E. Hoff (1960). Adaptive switching circuits. *Defense Technical Information Center*.

Zhang, T. (2004). Solving large scale linear prediction problems using stochastic gradient descent algorithms. In *Proceedings of the twenty-first international conference on Machine learning*, pp. 116. ACM.

A Appendix

A.1 Implicit stochastic approximation

Implicit stochastic approximation is defined as follows,

$$\theta_n^{\text{im}} = \theta_{n-1}^{\text{im}} - \gamma_n W_{\theta_n^{(*)}}, \text{ s.t.} \quad (2)$$

$$\mathbb{E}(\theta_n^{\text{im}} | \mathcal{F}_{n-1}) = \theta_{n-1}^{\text{im}} - \gamma_n h(\theta_n^{(*)}) = \theta_n^{(*)}, \quad (3)$$

where \mathcal{F}_{n-1} is the σ -algebra adapted to the sequence $\{\theta_0^{\text{im}}, \theta_1^{\text{im}}, \dots, \theta_{n-1}^{\text{im}}\}$.

A.2 Theory of implicit stochastic approximation

The symbol $\|\cdot\|$ denotes the L_1 vector/matrix norm. We also define the error random variables $\xi_n \triangleq W_{\theta_n^{(*)}} - h(\theta_{n-1}^{(*)})$, such that $\mathbb{E}(\xi_n | \mathcal{F}_{n-1}) = 0$. The parameter space for θ will be \mathbb{R}^p without loss of generality. For a positive scalar sequence a_n , the sequence $b_n = O(a_n)$ is such that $b_n \leq ca_n$, for some fixed $c > 0$, and every n ; the sequence $b_n = o(a_n)$ is such that $b_n/a_n \rightarrow 0$ in the limit where $n \rightarrow \infty$. $b_n \downarrow 0$ denotes a positive sequence decreasing towards zero. We further assume that implicit stochastic approximation (2) operates under a combination of the following assumptions.

Assumption 1. *It holds, $\gamma_n = \gamma_1 n^{-\gamma}$, $\gamma_1 > 0$ and $\gamma \in (1/2, 1]$.*

Assumption 2. *The regression function h is Lipschitz with parameter L , i.e., for all θ_1, θ_2 ,*

$$\|h(\theta_1) - h(\theta_2)\| \leq L \|\theta_1 - \theta_2\|.$$

Assumption 3. *Function h satisfies either*

(a) $(\theta - \theta_\star)^\top h(\theta) > 0$, for all θ , or

(b) $(\theta_n^{(*)} - \theta_\star)^\top h(\theta_n^{(*)}) \geq \delta_n \|\theta_n^{(*)} - \theta_\star\|^2$, where $\delta_n = \delta_1 n^{-\delta}$, $\delta_1 > 0$ and $0 < \gamma + \delta \leq 1$, for all n .

Assumption 4. *There exists a scalar potential $H : \mathbb{R}^p \rightarrow \mathbb{R}$ such that $\nabla H(\theta) = h(\theta)$, for all θ .*

Assumption 5. *There exists fixed $\sigma^2 > 0$ such that, for every n ,*

$$\mathbb{E}(\|\xi_n\|^2 | \mathcal{F}_{n-1}) \leq \sigma^2.$$

Assumption 6. *Let $\Xi_n \stackrel{\text{def}}{=} \mathbb{E}(\xi_n \xi_n^\top | \mathcal{F}_{n-1})$, then $\|\Xi_n - \Xi\| = O(1)$, and $\|\Xi_n - \Xi\| \rightarrow 0$ for fixed positive-definite matrix Ξ . Furthermore, if $\sigma_{n,s}^2 = \mathbb{E}(\mathbb{I}_{\|\xi_n\|^2 \geq s/\gamma_n} \|\xi_n\|^2)$, then for all $s > 0$, $\sum_{i=1}^n \sigma_{i,s}^2 = o(n)$ if $\gamma_n \propto n^{-1}$, or $\sigma_{n,s}^2 = o(1)$ otherwise.*

A.2.1 Convergence

Theorem 1. *Suppose that Assumptions 1, 2, 3(a), and 5 hold. Then the iterates θ_n^{im} of the implicit stochastic approximation procedure (2) converge almost-surely to θ_* ; i.e., $\theta_n^{\text{im}} \rightarrow \theta_*$, such that $h(\theta_*) = 0$, almost-surely.*

Proof. By definition (2) it follows

$$\|\theta_n^{\text{im}} - \theta_*\|^2 = \|\theta_{n-1}^{\text{im}} - \theta_*\|^2 - 2\gamma_n(\theta_{n-1}^{\text{im}} - \theta_*)^\top W_{\theta_n^{(*)}} + \gamma_n^2 \|W_{\theta_n^{(*)}}\|^2, \quad (4)$$

where $(\theta_{n-1}^{\text{im}} - \theta_n^{(*)}) = \gamma_n h(\theta_n^{(*)})$ by Eq.(3). Setting $(\theta_{n-1}^{\text{im}} - \theta_*) = (\theta_n^{(*)} - \theta_*) + (\theta_{n-1}^{\text{im}} - \theta_n^{(*)})$, yields

$$\begin{aligned} R_n &\triangleq \mathbb{E} \left((\theta_{n-1}^{\text{im}} - \theta_*)^\top W_{\theta_n^{(*)}} | \mathcal{F}_{n-1} \right) = (\theta_n^{(*)} - \theta_*)^\top h(\theta_n^{(*)}) + (\theta_{n-1}^{\text{im}} - \theta_n^{(*)})^\top h(\theta_n^{(*)}) \\ &= (\theta_n^{(*)} - \theta_*)^\top h(\theta_n^{(*)}) + \gamma_n \|h(\theta_n^{(*)})\|^2 > 0. \quad [\text{by Assumption 3(a)}] \end{aligned} \quad (5)$$

Through Eq.(3) we obtain

$$\begin{aligned} \|\theta_{n-1}^{\text{im}} - \theta_*\|^2 &= \|\theta_n^{(*)} - \theta_*\|^2 + 2\gamma_n h(\theta_n^{(*)})^\top (\theta_n^{(*)} - \theta_*) + \gamma_n^2 \|h(\theta_n^{(*)})\|^2, \\ &> \|\theta_n^{(*)} - \theta_*\|^2. \quad [\text{by Assumption 3(a)}] \end{aligned} \quad (6)$$

So, update (3) is non-expansive. Therefore,

$$\begin{aligned} \|h(\theta_n^{(*)})\| &= \|h(\theta_n^{(*)}) - h(\theta_*)\| \leq L \|\theta_n^{(*)} - \theta_*\| \quad [\text{by Assumption 2}] \\ &\leq L \|\theta_{n-1}^{\text{im}} - \theta_*\| \quad [\text{by non-expansiveness of update Eq.(3).}] \end{aligned} \quad (7)$$

Furthermore,

$$\begin{aligned} \mathbb{E} \left(\|W_{\theta_n^{(*)}}\|^2 | \mathcal{F}_{n-1} \right) &\stackrel{\text{def}}{=} \mathbb{E} \left(\|h(\theta_n^{(*)}) + \xi_n\|^2 | \mathcal{F}_{n-1} \right) \\ &= \|h(\theta_n^{(*)})\|^2 + \mathbb{E} \left(\|\xi_n\|^2 | \mathcal{F}_{n-1} \right) \\ &\leq L^2 \|\theta_{n-1}^{\text{im}} - \theta_*\|^2 + \sigma^2. \quad [\text{by Eq.(7) and Assumption 5}] \end{aligned} \quad (8)$$

Taking expectations in Eq.(4) conditional on \mathcal{F}_{n-1} and using Eq.(5) and Ineq.(8), we obtain

$$\mathbb{E} \left(\|\theta_n^{\text{im}} - \theta_*\|^2 | \mathcal{F}_{n-1} \right) \leq (1 + \gamma_n^2 L^2) \|\theta_{n-1}^{\text{im}} - \theta_*\|^2 - 2\gamma_n R_n + \gamma_n^2 \sigma^2. \quad (9)$$

We now use an argument —due to Gladyshev (1965)— that is also applicable to the classical Robbins-Monro procedure; see, for example, Benveniste et al. (1990, Section 5.2.2), or Ljung et al. (1992, Theorem 1.9). Random variable R_n is positive almost-surely by Ineq. (5), and $\sum \gamma_i = \infty$ and $\sum \gamma_i^2 < \infty$ by Assumption 1. Therefore, we can invoke the supermartingale lemma of Robbins and Siegmund (1985) to infer that $\|\theta_n^{\text{im}} - \theta_*\|^2 \rightarrow B > 0$ and $\sum \gamma_n R_n < \infty$, almost-surely. If $B \neq 0$ then $\liminf \|\theta_n^{\text{im}} - \theta_*\| > 0$, and thus the series $\sum_n \gamma_n R_n$ diverges by Ineq.(5) and $\sum \gamma_i = \infty$ (Assumption 1). Thus, $B = 0$. \square

A.3 Non-asymptotic analysis

Theorem 2. Suppose that Assumptions 1, 2, 3(a), 4, and 5 hold. Define $\Gamma^2 \triangleq \mathbb{E}(\|\theta_0^{\text{im}} - \theta_\star\|^2) + \sigma^2 \sum_{i=1}^\infty \gamma_i^2 + \gamma_1^2 \sigma^2$. Then, if $\gamma \in (2/3, 1]$, there exists $n_{0,1} < \infty$ such that, for all $n > n_{0,1}$,

$$\mathbb{E}(H(\theta_n^{\text{im}}) - H(\theta_\star)) \leq \left[\frac{2\Gamma^2}{\gamma\gamma_1} + o(1) \right] n^{-1+\gamma}.$$

If $\gamma \in (1/2, 2/3)$, there exists $n_{0,2} < \infty$ such that, for all $n > n_{0,2}$,

$$\mathbb{E}(H(\theta_n^{\text{im}}) - H(\theta_\star)) \leq \left[\Gamma\sigma\sqrt{L\gamma_1} + o(1) \right] n^{-\gamma/2}.$$

Otherwise, $\gamma = 2/3$ and there exists $n_{0,3} < \infty$ such that, for all $n > n_{0,3}$,

$$\mathbb{E}(H(\theta_n^{\text{im}}) - H(\theta_\star)) \leq \left[\frac{3 + \sqrt{9 + 4\gamma_1^3 L\sigma^2/\Gamma^2}}{2\gamma_1/\Gamma^2} + o(1) \right] n^{-1/3}.$$

Proof. By Eq.(3) and Assumption 3(a), $\theta_n^{(*)} + \gamma_n h(\theta_n^{(*)}) = \theta_{n-1}^{\text{im}}$ is equivalent to minimization $\theta_n^{(*)} = \arg \min_\theta \{ \frac{1}{2\gamma_n} \|\theta - \theta_{n-1}^{\text{im}}\|^2 + H(\theta) \}$. Therefore, comparing the values of the expression for $\theta = \theta_n^{(*)}$ and $\theta = \theta_{n-1}^{\text{im}}$, we obtain

$$H(\theta_n^{(*)}) + \frac{1}{2\gamma_n} \|\theta_n^{(*)} - \theta_{n-1}^{\text{im}}\|^2 \leq H(\theta_{n-1}^{\text{im}}). \quad (10)$$

Since $\theta_{n-1}^{\text{im}} - \theta_n^{(*)} = \gamma_n h(\theta_n^{(*)})$, Eq.(10) can be written as

$$H(\theta_{n-1}^{\text{im}}) - H(\theta_n^{(*)}) - \frac{1}{2}\gamma_n \|h(\theta_n^{(*)})\|^2 \geq 0. \quad (11)$$

We also have

$$\begin{aligned} H(\theta_n^{(*)}) - H(\theta_\star) &\leq h(\theta_n^{(*)})^\top (\theta_n^{(*)} - \theta_\star) \quad [\text{by Assumption 3(a)}] \\ H(\theta_n^{(*)}) - H(\theta_\star) &\leq \|h(\theta_n^{(*)})\| \cdot \|\theta_n^{(*)} - \theta_\star\| \quad [\text{by Cauchy-Schwartz}] \\ [\mathbb{E}(H(\theta_n^{(*)}) - H(\theta_\star))]^2 &\leq \mathbb{E}(\|h(\theta_n^{(*)})\|^2) \mathbb{E}(\|\theta_n^{(*)} - \theta_\star\|^2) \quad [\text{by Jensen's inequality}]. \end{aligned} \quad (12)$$

Furthermore,

$$\theta_n^{\text{im}} = \theta_{n-1}^{\text{im}} - \gamma_n (h(\theta_n^{(*)}) + \xi_n) = \theta_n^{(*)} - \gamma_n \xi_n. \quad [\text{by Eq.(3)}] \quad (13)$$

Therefore,

$$\begin{aligned} \mathbb{E}(\|\theta_n^{\text{im}} - \theta_\star\|^2) &= \mathbb{E}(\|\theta_n^{(*)} - \theta_\star\|^2) - 2\gamma_n \mathbb{E}((\theta_n^{(*)} - \theta_\star)^\top \xi_n) + \gamma_n^2 \mathbb{E}(\|\xi_n\|^2) \\ &= \mathbb{E}(\|\theta_n^{(*)} - \theta_\star\|^2) + \gamma_n^2 \mathbb{E}(\|\xi_n\|^2) \\ &\leq \mathbb{E}(\|\theta_{n-1}^{\text{im}} - \theta_\star\|^2) + \gamma_n^2 \sigma^2. \quad [\text{by Ineq.(6) and Assumption 5}] \\ &\leq \mathbb{E}(\|\theta_0^{\text{im}} - \theta_\star\|^2) + \sigma^2 \sum_{i=1}^n \gamma_i^2. \quad [\text{by induction.}] \end{aligned} \quad (14)$$

For notational convenience, define $h_n \stackrel{\text{def}}{=} \mathbb{E} (H(\theta_n^{\text{im}}) - H(\theta_\star))$ and $h_n^{(*)} \stackrel{\text{def}}{=} \mathbb{E} (H(\theta_n^{(*)}) - H(\theta_\star))$. Our goal is to derive a bound for h_n . Now, by Eq.(13) $\theta_n^{(*)} = \theta_n^{\text{im}} + \gamma_n \xi_n$. Since $\mathbb{E}(\xi_n | \mathcal{F}_{n-1}) = 0$, it follows from Assumption 5, $\mathbb{E}(\|\theta_n^{(*)} - \theta_\star\|^2) \leq \mathbb{E}(\|\theta_n^{\text{im}} - \theta_\star\|^2) + \gamma_n^2 \sigma^2$. Hence, using Ineq.(14),

$$\mathbb{E}(\|\theta_n^{(*)} - \theta_\star\|^2) \leq \mathbb{E}(\|\theta_0^{\text{im}} - \theta_\star\|^2) + \sigma^2 \sum_{i=1}^{\infty} \gamma_i^2 + \gamma_n^2 \sigma^2 \leq \Gamma^2, \quad (15)$$

by definition of Γ^2 . Furthermore, by convexity of H and Lipschitz continuity of h (Assumption 3(a)), and Assumption 5, we have

$$\begin{aligned} H(\theta_n^{\text{im}}) &= H(\theta_n^{(*)} - \gamma_n \xi_n) \\ H(\theta_n^{\text{im}}) &\leq H(\theta_n^{(*)}) - \gamma_n h(\theta_n^{(*)})^\top \xi_n + \gamma_n^2 \frac{L}{2} \|\xi_n\|^2 \quad [\text{by Lipschitz continuity}] \\ H(\theta_n^{\text{im}}) - H(\theta_\star) &\leq H(\theta_n^{(*)}) - H(\theta_\star) - \gamma_n h(\theta_n^{(*)})^\top \xi_n + \gamma_n^2 \frac{L}{2} \|\xi_n\|^2 \\ h_n &\leq h_n^{(*)} + \gamma_n^2 \frac{L\sigma^2}{2}. \quad [\text{by taking expectations.}] \end{aligned} \quad (16)$$

From Ineq.(12) and Ineq.(15),

$$\mathbb{E}(\|h(\theta_n^{(*)})\|^2) \geq \frac{1}{\Gamma^2} [\mathbb{E}(H(\theta_n^{(*)}) - H(\theta_\star))]^2 \stackrel{\text{def}}{=} \frac{1}{\Gamma^2} h_n^{(*)2}. \quad (17)$$

Now, in Ineq. (11), we subtract $H(\theta_\star)$ from the left-hand side, take expectations, and combine with (17) to obtain

$$h_{n-1} \geq h_n^{(*)} + \frac{1}{2\Gamma^2} \gamma_n h_n^{(*)2} \triangleq F_{\gamma_n}(h_n^{(*)}). \quad (18)$$

The function $F_{\gamma_n}(x)$ is monotone increasing because $h_n \geq 0$ and $h_n^{(*)} \geq 0$, since $H(\theta_\star)$ is minimum. Let $F_{\gamma_n}^{-1}$ denote its inverse, which is also monotone increasing. Thus $h_n^{(*)} \leq F_{\gamma_n}^{-1}(h_{n-1})$. Using Eq. (18) we can rewrite (16) as

$$h_n \leq F_{\gamma_n}^{-1}(h_{n-1}) + \gamma_n^2 \frac{L\sigma^2}{2}. \quad (19)$$

Ineq.(19) is our main recursion, since we want to upper-bound h_n . We will now try to find a base sequence $\{b_n\}$ such that $b_n \geq F_{\gamma_n}^{-1}(b_{n-1}) + \gamma_n^2 \frac{L\sigma^2}{2}$. Since one can take b_n to be increasing arbitrarily, we will try to find the smallest possible sequence $\{b_n\}$. To make our analysis more tractable we will search in the family of sequences $b_n = b_1 n^{-\beta}$, for various values $b_1, \beta > 0$. Then, b_n will be an upper-bound for h_n . To see this inductively, assume that $h_{n-1} \leq b_{n-1}$ and that h_n satisfies (19). Then, $h_n \leq F_{\gamma_n}^{-1}(h_{n-1}) + \gamma_n^2 \frac{L\sigma^2}{2} \leq F_{\gamma_n}^{-1}(b_{n-1}) + \gamma_n^2 \frac{L\sigma^2}{2} \leq b_n$, where the first inequality follows from the monotonicity of F_{γ_n} , and the second inequality follows from definition of b_n .

Now, the condition for b_n can be rewritten as $b_{n-1} \leq F_{\gamma_n}(b_n - \gamma_n^2 \frac{L\sigma^2}{2})$, and by definition of F_{γ_n} we get

$$b_{n-1} \leq b_n - \gamma_n^2 \frac{L\sigma^2}{2} + \gamma_n \frac{1}{2\Gamma^2} (b_n - \gamma_n^2 \frac{L\sigma^2}{2})^2 \quad (20)$$

Using $b_n = b_1 n^{-\beta}$ and $\gamma_n = \gamma_1 n^{-\gamma}$ (Assumption 1), we obtain

$$b_1[(n-1)^{-\beta} - n^{-\beta}] + \frac{L\sigma^2\gamma_1^2}{2} n^{-2\gamma} + \frac{L\sigma^2\gamma_1^3 b_1}{2\Gamma^2} n^{-\beta-3\gamma} - \frac{\gamma_1 b_1^2}{2\Gamma^2} n^{-2\beta-\gamma} - \frac{L^2\sigma^4\gamma_1^5}{8\Gamma^2} n^{-5\gamma} \leq 0. \quad (21)$$

We have $(n-1)^{-\beta} - n^{-\beta} < \frac{1}{1-\beta} n^{-1-\beta}$, for $n > 1$. Thus, it suffices to have

$$\frac{b_1}{1-\beta} n^{-1-\beta} + \frac{L\sigma^2\gamma_1^2}{2} n^{-2\gamma} + \frac{L\sigma^2\gamma_1^3 b_1}{2\Gamma^2} n^{-\beta-3\gamma} - \frac{\gamma_1 b_1^2}{2\Gamma^2} n^{-2\beta-\gamma} \leq 0, \quad (22)$$

where we dropped the $n^{-5\gamma}$ term without loss of generality. The positive terms in Ineq. (22) are $n^{-1-\beta}$, $n^{-2\gamma}$, and $n^{-\beta-3\gamma}$, and the only negative term is of order $n^{-2\beta-\gamma}$. In order to find the largest possible β to satisfy (22), one needs to equate the term $n^{-2\beta-\gamma}$ with the slowest possible term with a positive coefficient, i.e., set $2\beta + \gamma = \min\{1 + \beta, \beta + 3\gamma, 2\gamma\}$. However, $\beta + 3\gamma > 1 + \beta$ and $\beta + 3\gamma > 2\gamma$, and thus $2\beta + \gamma = \min\{1 + \beta, 2\gamma\}$, which implies only three cases:

- (a) $1 + \beta < 2\gamma$, and thus $2\beta + \gamma = 1 + \beta$, which implies $\beta = 1 - \gamma$. Also, $1 + \beta < 2\gamma \Rightarrow 2 - \gamma < 2\gamma$, and thus $\gamma \in (2/3, 1]$. In this case, b_1 will satisfy (22) for all $n > n_{0,1}$, for some $n_{0,1}$, if

$$\frac{b_1}{1-\beta} < \frac{\gamma_1 b_1^2}{2\Gamma^2} \Leftrightarrow b_1 > \frac{2\Gamma^2}{\gamma\gamma_1}. \quad (23)$$

- (b) $2\gamma < 1 + \beta$, and thus $2\beta + \gamma = 2\gamma$, which implies $\beta = \gamma/2$. Also, $1 + \beta > 2\gamma \Rightarrow 1 + \gamma/2 > 2\gamma$, and thus $\gamma \in (1/2, 2/3)$. In this case, b_1 will satisfy (22) for all $n > n_{0,2}$, for some $n_{0,2}$, if

$$\frac{\gamma_1^2 L\sigma^2}{2} < \frac{\gamma_1 b_1^2}{2\Gamma^2} \Leftrightarrow b_1 > \Gamma\sigma\sqrt{L\gamma_1}. \quad (24)$$

- (c) $2\gamma = 1 + \beta$, and thus $2\gamma = 1 + \beta = 2\beta + \gamma$, which implies $\gamma = 2/3$ and $\beta = 1/3$. In this case,

$$\frac{b_1}{1-\beta} + \frac{\gamma_1^2 L\sigma^2}{2} < \frac{\gamma_1 b_1^2}{2\Gamma^2}. \quad (25)$$

Because all constants are positive in Ineq.(25), including b_1 , it follows that

$$b_1 > \frac{3 + \sqrt{9 + 4\gamma_1^3 L\sigma^2/\Gamma^2}}{2\gamma_1/\Gamma^2}. \quad (26)$$

Remark. The constants $n_{0,1}, n_{0,2}, n_{0,3}$ depend on the problem parameters and the desired accuracy in the bounds of Theorem 2. Thus, it is straightforward to derive exact values for them. For example, consider case (a) and assume we picked b_1 such that $\frac{\gamma_1 b_1^2}{2\Gamma^2} - \frac{b_1}{1-\beta} = \epsilon > 0$. Ignoring the term $n^{-3\gamma-\beta}$ (for simplicity), Ineq.(22) becomes

$$\epsilon n^{-2+\gamma} \geq \frac{L\sigma^2\gamma_1^2}{2} n^{-2\gamma} \Rightarrow n \geq \left(\frac{L\sigma^2\gamma_1^2}{2\epsilon}\right)^c \equiv n_{0,1}, \quad (27)$$

where $c = 1/(3\gamma - 2) > 0$ since $\gamma \in (2/3, 1]$. If the desired accuracy is small, then ϵ is large and so the value $n_{0,1}$ will become smaller. Similarly, we can derive expressions for $n_{0,2}$ and $n_{0,3}$. \square

Theorem 3. Suppose that Assumptions 1, 3(b), and 5 hold, and define $\zeta_n \triangleq \mathbb{E}(\|\theta_n^{\text{im}} - \theta_\star\|^2)$ and $\kappa \triangleq 1 + 2\gamma_1\delta_1$. Then, if $\gamma + \delta < 1$, for every $n > 0$ it holds,

$$\zeta_n \leq e^{-\log \kappa \cdot n^{1-\gamma-\delta}} \zeta_0 + \sigma^2 \frac{\gamma_1 \kappa}{\delta_1} n^{-\gamma+\delta} + O(n^{-\gamma+\delta-1}).$$

Otherwise, if $\gamma = 1, \delta = 0$, it holds,

$$\zeta_n \leq e^{-\log \kappa \cdot \log n} \zeta_0 + \sigma^2 \frac{\gamma_1 \kappa}{\delta_1} n^{-1} + O(n^{-2}).$$

Proof. First we prove two lemmas that will be useful for Theorem 3.

Lemma 1. Consider a sequence b_n such that $b_n \downarrow 0$ and $\sum_{i=1}^\infty b_i = \infty$. Then, there exists a positive constant $K > 0$, such that

$$\prod_{i=1}^n \frac{1}{1+b_i} \leq \exp(-K \sum_{i=1}^n b_i). \quad (28)$$

Proof. The function $x \log(1 + 1/x)$ is increasing-concave in $(0, \infty)$. From $b_n \downarrow 0$ it follows that $\log(1 + b_n)/b_n$ is non-increasing. Consider the value $K \stackrel{\text{def}}{=} \log(1 + b_1)/b_1$. Then, $(1 + b_n)^{-1} \leq \exp(-K b_n)$. Successive applications of this inequality yields Ineq. (28). \square

Lemma 2. Consider sequences $a_n \downarrow 0, b_n \downarrow 0$, and $c_n \downarrow 0$ such that, $a_n = o(b_n)$, $\sum_{i=1}^\infty a_i \stackrel{\text{def}}{=} A < \infty$, and there is n' such that $c_n/b_n < 1$ for all $n > n'$. Define,

$$\delta_n \triangleq \frac{1}{a_n} (a_{n-1}/b_{n-1} - a_n/b_n) \text{ and } \zeta_n \triangleq \frac{c_n}{b_{n-1}} \frac{a_{n-1}}{a_n}, \quad (29)$$

and suppose that $\delta_n \downarrow 0$ and $\zeta_n \downarrow 0$. Pick a positive n_0 such that $\delta_n + \zeta_n < 1$ and $(1 + c_n)/(1 + b_n) < 1$, for all $n \geq n_0$.

Consider a positive sequence $y_n > 0$ that satisfies the recursive inequality,

$$y_n \leq \frac{1 + c_n}{1 + b_n} y_{n-1} + a_n. \quad (30)$$

Then, for every $n > 0$,

$$y_n \leq K_0 \frac{a_n}{b_n} + Q_1^n y_0 + Q_{n_0+1}^n (1 + c_1)^{n_0} A, \quad (31)$$

where $K_0 \stackrel{\text{def}}{=} (1 + b_1) (1 - \delta_{n_0} - \zeta_{n_0})^{-1}$, $Q_i^n = \prod_{j=i}^n (1 + c_j) / (1 + b_i)$, and $Q_i^n = 1$ if $n < i$, by definition.

Proof. We first consider two separate cases, namely, $n < n_0$ and $n \geq n_0$, and then we will combine the respective bounds.

Analysis for $n < n_0$. We first find a crude bound for Q_{i+1}^n . It holds,

$$Q_{i+1}^n \leq (1 + c_{i+1})(1 + c_{i+2}) \cdots (1 + c_n) \leq (1 + c_1)^{n_0}, \quad (32)$$

since $c_1 \geq c_n$ ($c_n \downarrow 0$ by definition) and there are no more than n_0 terms in the product. From Ineq. (30) we get

$$\begin{aligned} y_n &\leq Q_1^n y_0 + \sum_{i=1}^n Q_{i+1}^n a_i \quad [\text{by expanding recursive Ineq. (30)}] \\ &\leq Q_1^n y_0 + (1 + c_1)^{n_0} \sum_{i=1}^n a_i \quad [\text{using Ineq. (32)}] \\ &\leq Q_1^n y_0 + (1 + c_1)^{n_0} A. \end{aligned} \quad (33)$$

This inequality holds also for $n = n_0$.

Analysis for $n \geq n_0$. In this case, we have for all $n \geq n_0$,

$$\begin{aligned} (1 + b_1) (1 - \delta_n - \zeta_n)^{-1} &\leq K_0 \quad [\text{by definition of } n_0, K_0] \\ K_0 (\delta_n + \zeta_n) + 1 + b_1 &\leq K_0 \\ K_0 (\delta_n + \zeta_n) + 1 + b_n &\leq K_0 \quad [\text{because } b_n \leq b_1, \text{ since } b_n \downarrow 0] \\ \frac{1}{a_n} K_0 \left(\frac{a_{n-1}}{b_{n-1}} - \frac{a_n}{b_n} \right) + \frac{1}{a_n} K_0 \frac{c_n a_{n-1}}{b_{n-1}} + 1 + b_n &\leq K_0 \quad [\text{by definition of } \delta_n, \zeta_n] \\ a_n (1 + b_n) &\leq K_0 a_n - K_0 \left(\frac{(1 + c_n) a_{n-1}}{b_{n-1}} - \frac{a_n}{b_n} \right) \\ a_n &\leq K_0 \left(\frac{a_n}{b_n} - \frac{1 + c_n}{1 + b_n} \frac{a_{n-1}}{b_{n-1}} \right). \end{aligned} \quad (34)$$

Now we combine Ineqs. (34) and (30) to obtain

$$(y_n - K_0 \frac{a_n}{b_n}) \leq \frac{1 + c_n}{1 + b_n} (y_{n-1} - K_0 \frac{a_{n-1}}{b_{n-1}}). \quad (35)$$

For brevity, define $s_n \triangleq y_n - K_0 a_n / b_n$. Then, from Ineq. (35), $s_n \leq \frac{1+c_n}{1+b_n} s_{n-1}$, where $\frac{1+c_n}{1+b_n} < 1$ since $n \geq n_0$. Assume n_1 is the smallest integer such that $n_1 \geq n_0$ and $s_{n_1} \leq 0$ (existence of n_1 is not crucial.) For all $n \geq n_1$, it follows $s_n \leq 0$, and thus $y_n \leq K_0 a_n / b_n$

for all $n \geq n_1$. Alternatively, when $n_0 \leq n < n_1$, all s_n are positive. Using Ineq. (35) we have $s_n \leq (\prod_{i=n_0+1}^n \frac{1+c_i}{1+b_i}) s_{n_0} \stackrel{\text{def}}{=} Q_{n_0+1}^n s_{n_0}$, and thus

$$\begin{aligned} y_n - K_0 \frac{a_n}{b_n} &\leq Q_{n_0+1}^n s_{n_0} \quad [\text{by definition of } s_n] \\ y_n &\leq K_0 \frac{a_n}{b_n} + Q_{n_0+1}^n y_{n_0} \quad [\text{because } s_n \leq y_n] \\ y_n &\leq K_0 \frac{a_n}{b_n} + Q_1^n y_0 + Q_{n_0+1}^n (1+c_1)^{n_0} A. \quad [\text{by Ineq. (33) on } y_{n_0}]. \end{aligned} \quad (36)$$

Ineq. (36) holds for every n . \square

Corollary 1. In Lemma 2 assume $a_n = a_1 n^{-\alpha}$ and $b_n = b_1 n^{-\beta}$, and $c_n = 0$, where $\alpha > \beta$, and $a_1, b_1, \beta > 0$ and $1 < \alpha < 1 + \beta$. Then,

$$y_n \leq 2 \frac{a_1(1+b_1)}{b_1} n^{-\alpha+\beta} + \exp(-\log(1+b_1)n^{1-\beta})[y_0 + (1+b_1)^{n_0} A], \quad (37)$$

where $n_0 > 0$ and $A = \sum_i a_i < \infty$.

Proof. In this proof, we will assume, for simplicity, $(n-1)^{-c} - n^{-c} \leq n^{-1-c}$, $c \in (0, 1)$, for every $n > 0$. It is straightforward to derive an appropriate bound for each value of c . Furthermore, we assume $\sum_{i=1}^n i^{-\gamma} \geq n^{1-\gamma}$, for every $n > 0$. Formally, this holds for $n \geq n'$, where n' in practice is very small (e.g., $n' = 14$ if $\gamma = 0.1$, $n' = 5$ if $\gamma = 0.5$, and $n' = 9$ if $\gamma = 0.9$, etc.)

By definition,

$$\begin{aligned} \delta_n &\stackrel{\text{def}}{=} \frac{1}{a_n} \left(\frac{a_{n-1}}{b_{n-1}} - \frac{a_n}{b_n} \right) = \frac{1}{a_1 n^{-\alpha}} \frac{a_1}{b_1} ((n-1)^{-\alpha+\beta} - n^{-\alpha+\beta}) \\ &= \frac{1}{n^{-\alpha} b_1} [(n-1)^{-\alpha+\beta} - n^{-\alpha+\beta}] \\ &\leq \frac{1}{b_1} n^{-1+\beta}. \end{aligned} \quad (38)$$

Also, $\zeta_n = 0$ since $c_n = 0$. We can take $n_0 = \lceil (2/b_1)^{1/(1-\beta)} \rceil$, for which $\delta_{n_0} \leq 1/2$. Therefore, $K_0 \stackrel{\text{def}}{=} (1+b_1)(1-\delta_{n_0})^{-1} \leq 2(1+b_1)$; we can simply take $K_0 = 2(1+b_1)$. Since $c_n = 0$, $Q_i^n = \prod_{j=i}^n (1+b_j)^{-1}$. Thus,

$$\begin{aligned} Q_1^n &\geq (1+b_1)^{-n}, \text{ and} \\ Q_1^n &\leq \exp(-\log(1+b_1)/b_1 \sum_{i=1}^n b_i), \quad [\text{by Lemma 1.}] \\ Q_1^n &\leq \exp(-\log(1+b_1)n^{1-\beta}). \quad [\text{because } \sum_{i=1}^n i^{-\beta} \geq n^{1-\beta}.] \end{aligned} \quad (39)$$

Lemma 2 and Ineqs. (39) imply

$$\begin{aligned}
y_n &\leq K_0 \frac{a_n}{b_n} + Q_1^n y_0 + Q_{n_0+1}^n (1 + c_1)^{n_0} A \quad [\text{by Lemma 2}] \\
&\leq 2 \frac{a_1(1 + b_1)}{b_1} n^{-\alpha+\beta} + Q_1^n [y_0 + (1 + b_1)^{n_0} A] \quad [\text{by Ineqs. (39), } c_1 = 0] \\
&\leq 2 \frac{a_1(1 + b_1)}{b_1} n^{-\alpha+\beta} + \exp(-\log(1 + b_1)n^{1-\beta}) [y_0 + (1 + b_1)^{n_0} A],
\end{aligned} \tag{40}$$

where the last inequality also follows from Ineqs. (39). \square

Proof of Theorem 3. Now we are ready to prove the main theorem. By definition (2), $\theta_n^{\text{im}} = \theta_n^{(*)} - \gamma_n \xi_n$, and thus, by Assumption 5,

$$\mathbb{E} (||\theta_n^{\text{im}} - \theta_\star||^2) \leq \mathbb{E} (||\theta_n^{(*)} - \theta_\star||^2) + \gamma_n^2 \sigma^2 \tag{41}$$

By definition (3), $\gamma_n h(\theta_n^{(*)}) + \theta_n^{(*)} = \theta_{n-1}^{\text{im}}$, and thus

$$||\theta_{n-1}^{\text{im}} - \theta_\star||^2 = ||\theta_n^{(*)} - \theta_\star||^2 + 2\gamma_n (\theta_n^{(*)} - \theta_\star)^\top h(\theta_n^{(*)}) + \gamma_n^2 ||h(\theta_n^{(*)})||^2. \tag{42}$$

Therefore,

$$\begin{aligned}
||\theta_n^{(*)} - \theta_\star||^2 + 2\gamma_n (\theta_n^{(*)} - \theta_\star)^\top h(\theta_n^{(*)}) &\leq ||\theta_{n-1}^{\text{im}} - \theta_\star||^2 \\
||\theta_n^{(*)} - \theta_\star||^2 + 2\gamma_n \delta_n ||\theta_n^{(*)} - \theta_\star||^2 &\leq ||\theta_{n-1}^{\text{im}} - \theta_\star||^2 \quad [\text{by Assumption 3(b)}] \\
||\theta_n^{(*)} - \theta_\star||^2 &\leq \frac{1}{1 + 2\gamma_n \delta_n} ||\theta_{n-1}^{\text{im}} - \theta_\star||^2.
\end{aligned} \tag{43}$$

Combining Ineq.(41) and Ineq.(43) yields

$$\begin{aligned}
\mathbb{E} (||\theta_n^{\text{im}} - \theta_\star||^2) &= \mathbb{E} (||\theta_n^{(*)} - \theta_\star||^2) + \gamma_n^2 \sigma^2 \\
&\leq \frac{1}{1 + 2\gamma_n \delta_n} \mathbb{E} (||\theta_{n-1}^{\text{im}} - \theta_\star||^2) + \gamma_n^2 \sigma^2.
\end{aligned} \tag{44}$$

The final result of Theorem 3 is obtained through a direct application of Corollary 1 on recursion (44), by setting $y_n \equiv \mathbb{E} (||\theta_n^{\text{im}} - \theta_\star||^2)$, $b_n \equiv 2\gamma_n \delta_n$, and $a_n \equiv \gamma_n^2 \sigma^2$. The case where $\gamma = 1, \delta = 0$ only changes Ineq.(39) by replacing $\sum b_i$ with $\log n$. \square

A.4 Asymptotic distribution

Theorem 4. Suppose that Assumptions Assumption 1, Assumption 3(a), Assumption 5, and Assumption 6 hold. Then the iterate θ_n^{im} of implicit stochastic approximation (2) is asymptotically normal, such that

$$n^{\gamma/2} (\theta_n^{\text{im}} - \theta_\star) \rightarrow \mathcal{N}_p(0, \Sigma),$$

where Σ is the unique solution of

$$(\gamma_1 J_h(\theta_\star) - I/2)\Sigma + \Sigma(\gamma_1 J_h(\theta_\star)^\top - I/2) = \Xi, \tag{45}$$

where $J_h(\cdot)$ is the Jacobian of the regression function $h(\cdot)$, and I is the $p \times p$ identity matrix.

Proof. Convergence of $\theta_n^{\text{im}} \rightarrow \theta_\star$ is established from Theorem 1. By definition of the implicit stochastic approximation procedure (2),

$$\theta_n^{\text{im}} = \theta_{n-1}^{\text{im}} - \gamma_n(h(\theta_n^{(*)}) + \xi_n), \text{ and} \quad (46)$$

$$\theta_n^{(*)} + \gamma_n h(\theta_n^{(*)}) = \theta_{n-1}^{\text{im}}. \quad (47)$$

We use Eq. (47) and expand $h(\cdot)$ to obtain

$$\begin{aligned} h(\theta_n^{(*)}) &= h(\theta_{n-1}^{\text{im}}) - \gamma_n J_h(\theta_{n-1}^{\text{im}}) h(\theta_n^{(*)}) + \epsilon_n \\ h(\theta_n^{(*)}) &= (I + \gamma_n J_h(\theta_{n-1}^{\text{im}}))^{-1} h(\theta_{n-1}^{\text{im}}) + (I + \gamma_n J_h(\theta_{n-1}^{\text{im}}))^{-1} \epsilon_n, \end{aligned} \quad (48)$$

where $\|\epsilon_n\| = O(\gamma_n^2)$ by Theorem 3. By Lipschitz continuity of $h(\cdot)$ (Assumption 3(a)) and the almost-sure convergence of θ_n^{im} to θ_\star , it follows $h(\theta_{n-1}^{\text{im}}) = J_h(\theta_\star)(\theta_{n-1}^{\text{im}} - \theta_\star) + o(1)$, where $o(1)$ is a vector with vanishing norm. Therefore we can rewrite (48) as follows,

$$h(\theta_n^{(*)}) = A_n(\theta_{n-1}^{\text{im}} - \theta_\star) + O(\gamma_n^2), \quad (49)$$

such that $\|A_n - J_h(\theta_\star)\| \rightarrow 0$, and $O(\gamma_n^2)$ denotes a vector with norm $O(\gamma_n^2)$. Thus, we can rewrite (46) as

$$\theta_n^{\text{im}} - \theta_\star = (I - \gamma_n A_n)(\theta_{n-1}^{\text{im}} - \theta_\star) - \gamma_n \xi_n + O(\gamma_n^2). \quad (50)$$

The conditions for Fabian's theorem (Fabian, 1968, Theorem 1) are now satisfied, and thus $\theta_n^{\text{im}} - \theta_\star$ is asymptotically normal with mean zero, and variance that is given in the statement of Theorem 1 by Fabian (1968). \square