ASYMPTOTIC NORMALITY AND OPTIMALITY IN NONSMOOTH STOCHASTIC APPROXIMATION

By Damek Davis^{1, a}, Dmitriy Drusvyatskiy^{2, c} and Liwei Jiang^{1, b}

¹School of ORIE, Cornell University, ^adsd95@cornell.edu, ^blj282@cornell.edu

²Department of Mathematics, University of Washington, ^cddrusv@uw.edu

In their seminal work, Polyak and Juditsky showed that stochastic approximation algorithms for solving smooth equations enjoy a central limit theorem. Moreover, it has since been argued that the asymptotic covariance of the method is best possible among any estimation procedure in a local minimax sense of Hájek and Le Cam. A long-standing open question in this line of work is whether similar guarantees hold for important nonsmooth problems, such as stochastic nonlinear programming or stochastic variational inequalities. In this work, we show that this is indeed the case.

1. Introduction. Polyak and Juditsky [27] famously showed that the stochastic gradient method for minimizing smooth and strongly convex functions enjoys a central limit theorem: the error between the running average of the iterates and the minimizer, normalized by the square root of the iteration counter, converges to a normal random vector. Moreover, the asymptotic covariance matrix is in a precise sense "optimal" among any estimation procedure. A long standing open question is whether similar guarantees—asymptotic normality and optimality—exist for nonsmooth optimization and, more generally, for equilibrium problems. In this work, we obtain such guarantees under mild conditions that hold both in concrete circumstances (e.g., nonlinear programming) and under generic linear perturbations.

The types of problems we will consider are best modeled as stochastic variational inequalities. Setting the stage, consider the task of finding a solution X of the inclusion

$$(1.1) 0 \in \mathop{\mathsf{E}}_{z \sim P} A(x, z) + N_X(x).$$

Here, P is a probability distribution accessible only through sampling, $A(\cdot, z)$ is a smooth map for almost every $z \sim P$, and $N_X(x)$ denotes the normal cone to a closed set X. Stochastic variational inequalities (1.1) are ubiquitous in contemporary optimization. For example, optimality conditions for constrained optimization problems

$$\min_{X} \; \mathop{\mathsf{E}}_{z \sim P} \; f(x, z) \quad \text{subject to } x \in X,$$

fit into the framework (1.1) by setting $A(x, z) = \nabla f(x, z)$ in (1.1). More generally still, Nash equilibria $X = (x_1, \dots, x_n)$ of stochastic games are solutions of the system

$$X_j \in \underset{X_j \in X_j}{\operatorname{argmin}} \underset{z \sim P}{\mathsf{E}} f_j(x, z)$$
 for all $j = 1, \ldots, m$,

where f_j and X_j , respectively, are the loss function and the strategy set of player j. First order optimality conditions for these k coupled inclusions can be modeled as (1.1) by setting $[A(x, z)]_j := \nabla_{x_j} f_j(x, z)$ and $X := X_1, \ldots, X_m$.

Received January 2023; revised February 2024.

MSC2020 subject classifications. Primary 90C15; secondary 65K05, 65K10.

Key words and phrases. Stochastic gradient, asymptotic normality, local asymptotic minimax theory, variational inequality, active manifold.

There are two standard strategies for solving (1.1): sample average approximation (SAA) and the stochastic forward–backward algorithm (SFB). The former proceeds by drawing a batch of samples $Z_1, Z_2, \ldots, R^{-iid}P$ and finding a solution X_k to the empirical approximation

(1.2)
$$0 \in \frac{1}{k} \Big|_{i=1}^{k} A(x, z_i) + N_X(x).$$

In contrast, the stochastic forward–backward (SFB) algorithm proceeds in an online manner, drawing a single sample $Z_k \sim P$ in each iteration k and declaring the next iterate X_{k+1} as

$$(1.3) X_{k+1} \in P_X X_k - \alpha_k \cdot A(x_k, Z_k).$$

Here, $P_X(\cdot)$ denotes the nearest-point projection onto X. In the case of constrained optimization, $A(x, z) = \nabla f(x, z)$ is the gradient of some loss function f(x, z), and the process (1.3) reduces to the stochastic projected gradient algorithm. Online algorithms like SFB are usually preferable to SAA since each iteration is inexpensive and can be performed online, whereas SAA requires solving the auxiliary optimization problem (1.2). Although the asymptotic distribution of the SAA estimators is by now well understood [13, 14, 29], our understanding of the asymptotic performance of the SFB iterates is limited in nonsmooth and constrained settings. The goal of this paper is to fill this gap. The main result of our work is the following.

Under reasonable assumptions, the running average of the SFB iterates exhibits the same asymptotic distribution as SAA. Moreover, both SAA and SFB are asymptotically optimal in a locally minimax sense of Hájek and Le Cam [15, 31].

We next describe our results, and their consequences, in some detail. Namely, it is classically known (e.g., [13, 14, 29]) that the asymptotic performance of SAA (1.2) is strongly influenced by the sensitivity of the solution X to perturbations of the left-hand side of (1.1). In order to isolate this effect, let S(V) consist of solutions X to the perturbed system

$$v \in \underset{z \sim P}{\mathbb{E}} A(x, z) + N_X(x).$$

Throughout, we will assume that the solutions S(v) vary smoothly near X. More precisely, we will assume that the graph of S locally around (0, X) coincides with the graph of some smooth map $\sigma(\cdot)$. In the language of variational analysis [9], the map $\sigma(\cdot)$ is called a smooth localization of S around (0, X). It is known that this assumption holds in a variety of concrete circumstances and under generic linear perturbations of semialgebraic problems [10].

Let us next provide the context and state our results. It is known from [14, 29] that under mild assumptions, the solutions X_k of SAA (1.2) are asymptotically normal:

(1.4)
$$\sqrt{k} x_k - x \xrightarrow{D} N 0, \nabla \sigma (0) \cdot Cov A x, z \cdot \nabla \sigma (0) .$$

Thus the Jacobian of the solution map $\nabla \sigma$ (0) appears in the asymptotic covariance of the SAA estimator. In fact, we will argue that this is unavoidable. The first contributions of our work is that we prove that the asymptotic performance of SAA is locally minimax optimal—in the sense of Hájek and Le Cam [15, 31]—among all estimation procedures. Roughly speaking, this means that for any estimation procedure that outputs χ_k based on χ_k samples, there exists a sequence of perturbations χ_k with $\frac{dP_k}{dP} = 1 + O(k^{-1/2})$, such that the performance of χ_k on the perturbed sequence of problems is asymptotically no better than the performance of SAA on the perturbed problems. We note that the analogous lower bound for stochastic nonlinear programming was obtained earlier in [12], and our arguments are motivated by the techniques therein. The fact that the SFB algorithm for smooth problems is asymptotically optimal was proved in [4], Theorem 5.6, by verifying that SFB is asymptotically equivariant

in law; we follow a similar argument here. Aside from the lower bound, the main result of our work is to show that under reasonable assumptions, the running average of the SFB iterates enjoys the same asymptotics as (1.4) and is thus asymptotically optimal.

The guarantees we develop are already interesting for stochastic nonlinear programming:

(1.5)
$$\min_{X} f(x) = \mathop{\mathbb{E}}_{z \sim P} f(x, z) \quad \text{subject to} \quad g_i(x) \le 0 \quad \forall i = 1, \dots, m.$$

Here each g_i is a smooth function and the map $x \to f(x, z)$ is smooth for a.e. $z \sim P$. The optimality conditions for this problem can be modeled as the variational inequality (1.1) under the identification $A(x, z) = \nabla f(x, z)$ and $X = \{x : g_i(x) \le 0 \ \forall i = 1, \ldots, m\}$ The stochastic forward–backward algorithm then becomes the stochastic projected gradient method. Our results imply that under the three standard conditions—linear independence of active gradients, strict complementarity, and strong second-order sufficiency—the running average of the SFB iterates $\overline{X}_k = \frac{1}{k} \int_{i=1}^{k} X_i$ is asymptotically normal and optimal:

$$\sqrt[4]{\overline{k}} \ \overline{X}_k - \chi \quad \stackrel{D}{\longrightarrow} N \ 0, \ \nabla \sigma \ (0) \cdot \operatorname{Cov} \ \nabla f \quad X \ , \ Z \ \cdot \nabla \sigma \ (0) \ .$$

Moreover, as is classically known, the Jacobian $\nabla \sigma$ (0) admits an explicit description as

$$\nabla \sigma (0) = P_T \nabla_{xx}^2 L \ x \ , y \ P_T \ ^\dagger,$$

where $\nabla_{xx}^2 L(x,y)$ is the Hessian of the Lagrangian function, the symbol † denotes the Moore–Penrose pseudoinverse, and P_T is the projection onto the linear subspace $\{\nabla g_i(x)\}_{i\in I}^\perp$ and $I=\{i:g_i(x)=0\}$ is the set of active indices. An illustrative example of the announced result is depicted in Figure 1, which plots the performance of the projected stochastic gradient method for minimizing a linear function over the intersection of two balls. A further illustration for a nonconvex problem of sparse recovery is depicted in Figure 2. This result may be surprising in light of the existing literature. Namely, Duchi and Ruan [12] uncover a striking gap between the estimation quality of SAA and at least one standard online method, called dual averaging [24, 33], for stochastic nonlinear optimization. Indeed, even for the problem of minimizing the expectation of a linear function over a ball, the dual averaging method exhibits a suboptimal asymptotic covariance ([12], Section 5.2). ¹ In contrast, we see that the stochastic projected gradient method is asymptotically optimal.

Let us now return to the general problem (1.1) and the stochastic forward–backward algorithm (1.3). In order to derive the claimed asymptotic guarantees for SFB, we will impose a few extra assumptions. First, in addition to assuming that $\sigma(\cdot)$ is smooth near the origin, we will assume that there exists a neighborhood U of the origin such that $\sigma(U)$ is a smooth manifold. This assumption is mild, since it holds automatically for example, if the matrix $\nabla \sigma(\cdot)$ has constant rank on a neighborhood of the origin. In the language of [11], the set $M = \sigma(U)$ is called an *active manifold* around \overline{X} . Returning to the case of stochastic nonlinear programming, the active manifold is simply the zero-set of the active inequalities

$$M = x : g_i(x) = 0 \ \forall i \in I .$$

See Figure 1 for an illustration. Variants of active manifolds have been extensively studied in nonlinear programming, under the names of identifiable surfaces [32], partly smooth sets [19], UV-structures [18, 21], $g \circ F$ decomposable functions [30], and minimal identifiable sets [11].

The main idea of our argument is to relate the nonsmooth dynamics of SFB to a smooth stochastic approximation algorithm on M. More precisely, we will show that under mild

 $^{^{1}}$ In contrast, in the special case that X is polyhedral and convex, the dual averaging method is optimal [12].

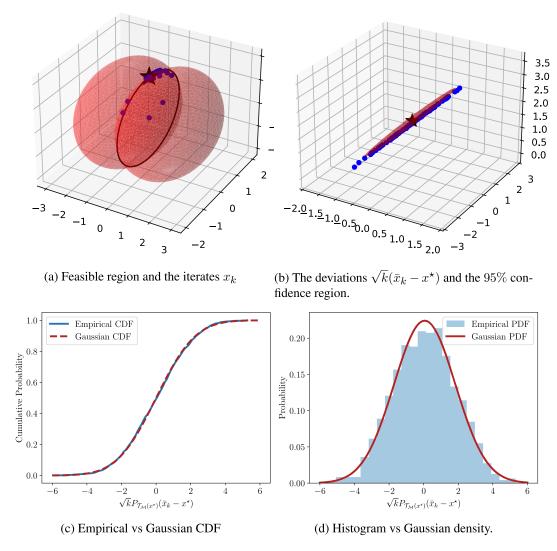
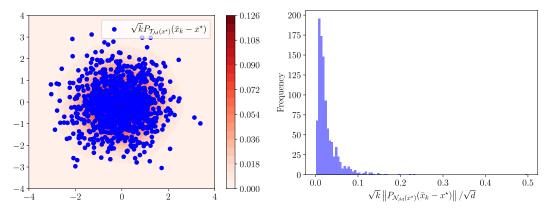


Fig. 1. The stochastic projected gradient method for minimizing Eg[$-x_1 + g$, x] over the intersection of two balls centered around(-1, 0, 0) and (1, 0, 0) of radius two. The expectation is taken over a Gaussiag $\sim N(0, 1)$. The optimal solution $(0, 0, \overline{3})$ (marked with a star) lies on the active manifold M, which is a circle depicted in black. The figure on the top left depicts the iterates generated by a single run of the process initialized at the origin with stepsize $n_k = k^{-3/4}$ and executed for 1000 iterations. The figure on the top right depicts the rescaled deviations $\overline{k}(\overline{x}_k - x)$ taken over 100 runs with $K = 10^6$. The two figures clearly show that the iterates rapidly approach the active manifold and asymptotically the deviations $\overline{k}(\overline{x}_k - x)$ are supported only along the tangent space to M at X. The two figures on the second row show the histogram and the empirical CDF, respectively, of the tangent components $\overline{k}P_{T_M}(x)(\overline{x}_k - x)$, overlaid with the analogous functions for a Gaussian.

conditions, the shadow sequence $\mathcal{Y}_k := P_M(X_k)$ along the manifold M behaves smoothly up to a small error

(1.7)
$$y_{k+1} = y_k - \alpha_k P_{T_M(y_k)} A(y_k, Z_k) + o(\alpha_k),$$

where T_M (Y_k) denotes the tangent space of M at Y_k . Consequently, we may build on the techniques of Polyak and Juditsky [27] to obtain the asymptotics of the shadow sequence Y_k , and then infer information about the original iterates X_k . We note that in the constrained optimization setting, the iteration (1.7) becomes an inexact Riemannian gradient method on the restriction of Y_k to Y_k .



- (a) Kernel density estimation on tangent deviations.
- (b) Histogram of normal deviations.

Fig. 2. The stochastic projected gradient—method for minimizing $E(a,b)[(a,x-b)^2]$ over the $_0$ ball $X = \{x : x \mid_{0} \le 2\}$. Here $a \sim N(0,1)$ and b = a,x+g where $g \sim N(0,1)$ and $x := e_1 + e_2$, the sum of the first two standard basis vectors; in this example, d = 20. The optimal solution x lies on the active manifold $M = \text{span}\{e_1, e_2\}$. The figure on the left—depicts a kernel—density estimation of—the rescaled deviations $\overline{K} \cdot PT_M(x)(\overline{X}K - x)$ —taken over 1000 runs of—SGD (Gaussian kernel, bandwidth 0.5); here, the method is initialized at the origin with stepsize $T_1 = x^{-3/4}$ and ran for $T_2 = x^{-1/4}$ and ran for $T_3 = x^{-1/4}$ and ran for

The validity of (1.7) relies on two extra conditions, introduced in [6], which relate the "first-order" behavior of X to that of M. Namely, for $x \in X$ and $y \in M$ near X, we assume:

•
$$N \times (y) \cap \mathbb{S}^{d-1}, x - y = o(x - y)$$

• $N \times (x) \cap \mathbb{S}^{d-1} \subseteq N_M(y) + O(x - y)B$

[(b)-regularity]
[strong (a)-regularity]

Here S^{d-1} and B are the unit sphere and the closed unit ball in R^d , respectively. The (b)-regularity condition simply asserts that the secant line joining $x \in X$ and $y \in M$ becomes tangent to M as X and Y tend to the same point near X. The strong-(a) regularity in contrast asserts that the normal cone $N_X(x)$ is contained in $N_M(y)$ up to a linear error O(x-y)—a kind of Lipschitz condition. The two regularity conditions are introduced and thoroughly developed in [6], with numerous examples and calculus rules presented. In particular, both conditions hold automatically for stochastic nonlinear programming.

- 1.1. *Outline*. The outline of the rest of the paper is as follows. Section 2 presents the basic notation and constructions that will be used in the paper. Existence of smooth localizations $\sigma(\cdot)$ is a central assumption of our work. Section 3 develops asymptotic convergence guarantees for SAA, which motivate much of the subsequent sections. Section 4 presents the classes of algorithms that we consider. Section 5 states the main result on asymptotic normality of iterative methods. Section 6 present shows that SAA and SFB are both asymptotically local minimax optimal in the sense of Hájek and Le Cam.
- **2. Notation and basic constructions.** This section records basic notation that we will use throughout the paper. To this end, the symbo \mathbb{R}^d will denote a Euclidean space with inner product \cdot , \cdot and the induced norm $x = \overline{X, X}$. The symbol B will stand for the closed unit ball in \mathbb{R}^d , while $B_r(X)$ will denote the closed ball of radius f around a point f. For any function $f: \mathbb{R}^d \to \mathbb{R} \cup \{+\infty\}$, the *domain*, *graph*, and *epigraph* are defined as

$$\operatorname{dom} f := \quad x \in \operatorname{R}^d : f(x) < \infty \quad ,$$

$$gph f := X, f(x) \in \mathbb{R}^d \times \mathbb{R} : x \in dom^f$$
,
 $epi f := (x, r) \in \mathbb{R}^d \times \mathbb{R} : r \ge f(x)$,

respectively. We say that f is *closed* if epi f is a closed set, or equivalently if f is lower-semicontinuous. The *proximal map* of f with parameter $\alpha > 0$ is given by

$$\operatorname{prox}_{\alpha f}(x) := \underset{y}{\operatorname{argmin}} f(y) + \frac{1}{2\alpha}y - x^{-2}$$
.

The *distance* and the *projection* of a point $x \in \mathbb{R}^d$ onto a set $Q \subseteq \mathbb{R}^d$ are

$$d(x, Q) := \inf_{y \in Q} y - x$$
 and $P_Q(x) := \underset{v \in O}{\operatorname{argmin}} y - x$,

respectively. The indicator function of Q, denoted by $\delta_Q(\cdot)$, is defined to be zero on Q and $+\infty$ off it. The symbol O(h) stands for any function $O(\cdot)$ satisfying $O(h)/h \to 0$ as $h \to 0$.

2.1. Smooth manifolds. Next, we recall a few definitions from smooth manifolds; we refer the reader to [2, 16] for details. Throughout the paper, all smooth manifolds M are assumed to be embedded in \mathbb{R}^d and we consider the tangent and normal spaces to M as subspaces of \mathbb{R}^d . Thus, a set $M \subset \mathbb{R}^d$ is a C^p manifold (with $p \ge 1$) if around any point $x \in M$ there exists an open neighborhood $U \subset \mathbb{R}^d$ and a C^p -smooth map F from U to some Euclidean space \mathbb{R}^n such that the Jacobian $\nabla F(x)$ is surjective and equality $M \cap U = F^{-1}(0)$ holds. Then F = 0 are called the local defining equations for M, and the tangent and normal spaces to M at X are defined by $T_M(x) := \text{Null}(\nabla F(x))$ and $N_M(x) := (T_M(x))^{\perp}$, respectively. Note that for C^p manifolds M with $p \ge 1$, the projection P_M is C^{p-1} -smooth on a neighborhood of each point X in M and is C^p -smooth on a neighborhood of the origin in the tangent space $T_M(x)$ [22]. Moreover, the inclusion range $(\nabla P_M(x)) \subseteq T_M(P_M(x))$ holds for all X near M and $\nabla P_M(x) = P_{T_M(x)}$ holds for all $X \in M$ ([20], Lemma 2.1). Let $M \subset \mathbb{R}^d$ be a C^p -manifold for some $p \ge 1$. Then a function $f : M \to \mathbb{R}$ is called

Let $M \subset \mathbb{R}^d$ be a C^p -manifold for some $p \geq 1$. Then a function $f \colon M \to \mathbb{R}$ is called C^p -smooth around a point $x \in M$ if there exists a C^p function $f \colon U \to \mathbb{R}$ defined on an open neighborhood $U \subset \mathbb{R}^d$ of X and that agrees with f on $U \cap M$. Then the *covariant gradient of* f at X is the vector $\nabla_M f(x) := P \ T_M(x) (\nabla f(x))$. When f and f are f are f and f are f and f are f and f are f are f and f are f are f and f are f are f and f are f are f and f are f and f are f are f and f are f and f are f are f and f are f and f are f are f and f are f are f and f are f and f are f and f are f and f are f are f and f are f and f are f are f and f are f are f and f are f and f are f are f and f are f are f and f are f and f are f are f are f and f are f and f are f are f and f are f and f are f are f and f are f are f are f and f are f are f are f are f and f are f are f and f are f and f are f are f are f are f are f are f and f are f are f are f are f and f are f ar

$$\frac{d^2}{dt^2}f P_M(x+tu) = \nabla_M^2 f(x)u, u \text{ for all } u \in T_M(x).$$

If M is \mathbb{C}^3 -smooth, then we can write $\nabla^2_M f(x)$ simply as

$$\nabla_{M}^{2} f(x) = P \tau_{M}(x) \nabla^{2} (f \circ P_{M})(x) P \tau_{M}(x),$$

while regarding the right-hand side as a linear operator on $T_M(x)$.

A map $F: M \to \mathbb{R}^m$ is called C^p smooth near a point X if there exists a map $\hat{F}: U \to \mathbb{R}^d$ defined on some neighborhood $U \subset \mathbb{R}^d$ of X that agrees with F on M near X. In this case, we define the *covariant Jacobian* $\nabla F(X): T_M(X) \to \mathbb{R}^m$ by the expression $\nabla_M F(X)u = \nabla F(X)u$ for all $u \in T_M(X)$. An easy computation shows that in the particular case when $F(X) = \nabla_M f(X)$ for a C^3 -smooth function $f: M \to \mathbb{R}$, the quadratic form defined by $\nabla_M F(X)$ on $T_M(X)$ coincides with $\nabla_M^2 f(X)$. More precisely, equality holds:

$$P_{T_M(x)} \cdot \nabla_M F(x) \cdot P_{T_M(x)} = P_{T_M(x)} \cdot \nabla_M^2 f(x) \cdot P_{T_M(x)}$$

2.2. *Normal cones and subdifferentials.* Next, we require a few basic constructions of nonsmooth and variational analysis. Our discussion will follow mostly closely Rockafellar—Wets [28]. Other influential treatments of the subject include [1, 3, 23, 25]. The *Fréchet normal cone* to a set $Q \subseteq \mathbb{R}^d$ at a point $x \in \mathbb{R}^d$, denoted $\hat{N}_Q(x)$, consists of all vectors $v \in \mathbb{R}^d$ satisfying

$$(2.1) v, y - x \le 0 y - x as y \to x in Q.$$

Thus V lies in $\hat{N}_Q(X)$ if up to first-order V makes an obtuse angle with all directions pointing from X into Q. Generally, Fréchet normals are highly discontinuous with respect to perturbations of the basepoint X. Consequently, we enlarge the Fréchet normal cone as follows. The *limiting normal cone* to Q at $X \in Q$, denoted by $N_Q(X)$, consists of all vectors $V \in \mathbb{R}^d$ for which there exist sequences $X_i \in Q$ and $V_i \in \hat{N}_Q(X_i)$ satisfying $(X_i, V_i) \to (X, V)$.

The analogous of normal cones for functions are subdifferentials, or sets of generalized derivatives. Namely, consider a function $f: \mathbb{R}^d \to \mathbb{R} \cup \{\infty\}$ and a point $x \in \text{dom}^f$. The *Fréchet subdifferential of* f at X, denoted $\partial f(x)$, consists of all vectors $v \in \mathbb{R}^d$ satisfying the approximation property:

$$f(y) \ge f(x) + v, y - x + o \qquad y - x$$
 as $y \to x$.

The *limiting subdifferential of* f *at* X, denoted $\partial f(X)$, consists of all $V \in \mathbb{R}^d$ such that there exist sequences $X_i \in \mathbb{R}^d$ and Fréchet subgradients $V_i \in \partial f(X_i)$ satisfying $(X_i, f(X_i), V_i) \to (X, f(X), V)$ as $i \to \infty$. A point X satisfying $0 \in \partial f(X)$ is called *critical* for f. The primary goal of algorithms for nonsmooth optimization is the search for critical points.

2.3. *Active manifolds*. Critical points of typical nonsmooth functions lie on a certain manifold that captures the activity of the problem in the sense that critical points of slight linear perturbation of the function do not leave the manifold. Such active manifolds have been modeled in a variety of ways, including identifiable surfaces [32], partial smoothness [19], UV-structures [18, 21], $g \circ F$ decomposable functions [30], and minimal identifiable sets [11]. In this work, we adopt the following formal model of activity, explicitly used in [11].

DEFINITION 2.1 (Active manifold). Consider a function $f: \mathbb{R}^d \to \mathbb{R} \cup \{\infty\}$ and fix a set $M \subseteq \text{dom } f$ containing a point \overline{X} satisfying $0 \in \partial f(\overline{X})$. Then M is called an *active* C^p -manifold around \overline{X} if there exists a constant 0 satisfying the following:

- (smoothness) The set M is a C^p manifold near \overline{X} and the restriction of f to M is C^p -smooth near \overline{X} .
- (sharpness) The lower bound holds:

$$\inf \ v: v \in \partial f(x), \ x \in U \setminus M \ > 0,$$

where we set $U = \{x \in B \ (x\bar{)} : |f(x) - f(\bar{x})| < \}$

More generally, we say that M is an active manifold for f at \overline{X} for $\overline{v} \in \partial f(\overline{X})$ if M is an active manifold for the tilted function $f_{v}(x) = f(x) - v$, x at \overline{X} .

The sharpness condition simply means that the subgradients of f must be uniformly bounded away from zero at points off the manifold that are sufficiently close to \overline{X} in distance and in function value. The localization in function value can be omitted for example if f is weakly convex or if f is continuous on its domain; see [11] for details.

Intuitively, the active manifold has the distinctive feature that the function varies smoothly along the manifold and grows linearly normal to it; see Figure 3 for an illustration. This is summarized by the following theorem from [5], Theorem D.2.

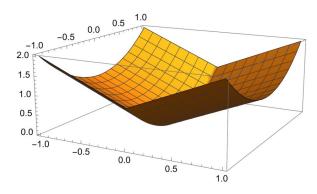


Fig. 3. $f(x_1, x_2) = |x_1| + x_2^2$.

PROPOSITION 2.2 (Identification implies sharpness). Suppose that a closed function $f: \mathbb{R}^d \to \mathbb{R} \cup \{\infty\}$ admits an active manifold M at a point \overline{X} satisfying $0 \in \hat{\partial} f(\overline{X})$. Then there exist constants C, 0 such that

$$(2.2) f(x) - f P_M(x) \ge c \cdot \operatorname{dist}(x, M) \quad \forall x \in B(\bar{x}).$$

Notice that there is a nontrivial assumption $0 \in \hat{\partial} f(\overline{X})$ at play in Proposition 2.2. Indeed, under the weaker inclusion $0 \in \partial f(\overline{X})$ the growth condition (2.2) may easily fail, as the univariate example f(x) = -|x| shows. It is worthwhile to note that under the assumption $0 \in \hat{\partial} f(\overline{X})$, the active manifold is locally unique around \overline{X} ([11], Proposition 8.2).

In order to make progress, we will require two extra conditions to hold along the active manifold that tightly couple the subgradients of f on and off the manifold. Although these two conditions, introduced in [6], may look formidable, they are very mild indeed.

The motivation for the first regularity condition stems from a weakening of Taylor's theorem that is appropriate for nonsmooth functions. Namely, recall that any C^1 -smooth function f satisfies the first-order approximation property

(2.3)
$$f(y) = f(x) + \nabla f(x), y - x + o y - x$$
 as $x, y \to \overline{x}$.

This estimate is fundamental to optimization theory and practice since it quantifies the approximation qualify of the linear model of f furnished by the gradient. When f is nonsmooth, the analogue of (2.3) with subgradients replacing the gradient can not possibly hold uniformly over all points X and Y near X, since it would imply differentiability. Instead, a reasonable requirement is to only require (2.3) to hold with points Y lying on the active manifold M. In fact, for our purposes, it suffices to replace the equality with an inequality.

DEFINITION 2.3 ((b_{\leq}) -regularity). Consider a function $f: \mathbb{R}^d \to \mathbb{R} \cup \{\infty\}$ that is locally Lipschitz continuous on its domain. Fix a set $M \subset \text{dom } f$ that is a C^1 manifold around \overline{X} and such that the restriction of f to M is C^1 -smooth near \overline{X} . We say that f is (b_{\leq}) -regular along M at \overline{X} if there exists 0 such that the estimates

(2.4)
$$f(y) \ge f(x) + v, y - x + 1 + v \cdot o y - x$$

hold for all $x \in \text{dom } f \cap B$ $(\vec{x}), y \in M \cap B$ $(\vec{x}), \text{ and } v \in \partial f(x)$.

Importantly, condition (b_{\leq}) along an active manifold M directly implies that all negative subgradients of f, taken at points X near \overline{X} , point towards the active manifold M. Indeed, this is a direct consequence of Proposition 2.2, and is summarized in the following corollary. Consequently, subgradient type methods always move in the direction of the active manifold—clearly a desirable property.

COROLLARY 2.4 (Proximal aiming). Consider a closed function $f: \mathbb{R}^d \to \mathbb{R} \cup \{\infty\}$ that admits an active C^1 -manifold M at a point \overline{X} satisfying $0 \in \widehat{\partial} f(\overline{X})$. Suppose that f is locally Lipschitz continuous on its domain and that f is (b_{\leq}) -regular along M near \overline{X} . Then, there exists a constant $\mu \geq 0$ such that, the estimate

(2.5)
$$V, X - P_M(X) \ge \mu \cdot \text{dist}(X, M) + 1 + V \cdot o \text{dist}(X, M)$$
,

holds for all $x \in \text{dom}^f$ near \overline{X} and for all $v \in \partial f(x)$.

The second regularity condition has a different flavor, stemming from a weakening of Lipschitz continuity of the gradient. Namely, nonsmoothness by its nature implies that the deviation $\partial f(x) - \partial f(y)$ is not controlled well by the distance in the arguments x - y. On the other hand, it turns out that if we take $y \in M$ and only look at the subgradient deviations in tangent directions, the error $P_{T_M}(y)[\partial f(x) - \partial f(y)]$ is typically linearly bounded by x - y. Moreover, in typical circumstances $P_{T_M}(y)[\partial f(y)]$ consists only of a single vector, the covariant gradient $\nabla_M f(y)$.

DEFINITION 2.5 (Strong (a)-regularity). Consider a function $f: \mathbb{R}^d \to \mathbb{R} \cup \{\infty\}$ that is locally Lipschitz continuous on its domain. Fix a set $M \subset \text{dom } f$ that is a C^1 manifold around \overline{X} and such that the restriction of f to M is C^1 -smooth near \overline{X} . We say that f is strongly (a)-regular along M near \overline{X} if there exist constants C, $0 \ge 0$

$$(2.6) P_{T_M(y)} v - \nabla_M f(y) \le C 1 + v \quad x - y,$$

for all $x \in \text{dom } f \cap B$ $(\vec{X}), y \in M \cap B$ $(\vec{X}), \text{ and } v \in \partial f(x)$.

The two regularity conditions easy extend to sets through their indicator functions. Namely, we say that a set $Q \subset \mathbb{R}^d$ is (b_{\leq}) -regular (respectively strongly (a)-regular) along a C^1 manifold $M \subset Q$ at $\overline{x} \in M$ if the indicator function δ_Q is (b_{\leq}) -regular (respectively, strongly (a)-regular) along M at \overline{x} .

The paper [6] presents a wide array of functions that admit active manifolds along which both conditions (b_{\leq}) and strong (a) hold. Here, we discuss in detail a single example of nonlinear programming, and refer the reader to [6] for many more examples.

EXAMPLE 1 (Nonlinear programming). Consider the problem of nonlinear programming

(2.7)
$$\min_{X} f(x)$$
$$s \cdot t \cdot g_i(x) \le 0 \quad \text{for } i = 1, \dots, m,$$
$$g_i(x) = 0 \quad \text{for } i = m + 1, \dots, n,$$

where f and g_i are C^p -smooth functions on R^d . Let X denotes the set of all feasible points to the problem. Consider now a point $\overline{x} \in X$ that is critical for the function $f + \delta x$ and define the active index set

$$I = i : g_i(\overline{x}) = 0.$$

Suppose the following is true:

• (*LICQ*) the gradients $\{\nabla g_i(\bar{x})\}_{i\in I}$ are linearly independent.

Then the set

$$M = x : q_i(x) = 0 \ \forall i \in I$$

is a C^p smooth manifold locally around \bar{X} . Moreover, all three functions f, δ_X , and $f + \delta_X$ are (b_{\leq}) -regular and strongly (a)-regular along M near \bar{X} . In order to ensure that M is an active manifold of $f + \delta x$, an extra condition is required. Define the Lagrangian function

$$L(x, y) := f(x) + \int_{i=1}^{n+m} y_i g_i(x).$$

Criticality of \overline{X} and LICQ ensures that there exists a (unique) Lagrange multiplier vector $\overline{y} \in \mathbb{R}^m_+ \times \mathbb{R}^n$ satisfying $\nabla_X L(\overline{x}, \overline{y}) = 0$ and $\overline{y}_i = 0$ for all $i \not \in I$. Suppose the following standard assumption is true:

• (Strict complementarity) $\overline{y_i} > 0$ for all $i \in I \cap \{1, ..., m\}$

Then M is indeed an active C^p manifold for $f + \delta x$ at \overline{x} .

EXAMPLE 2 (1-regularization). Consider the stochastic optimization problem with regularization

(2.8)
$$\min g(x) = f(x) + \lambda x$$
 1'

where $f(x) = E_{z \in P}[f(x, z)]$ is a C^p -smooth function in \mathbb{R}^d . Consider now a point $\overline{x} \in \mathbb{R}^d$ that is critical for the function g and define the index set f(x) = f(x). Then the set

$$M = \{x : x_i = \overline{X}_i, \forall i \in I \}$$

is an affine space, hence a smooth manifold. Note that the definition of criticality ensures that $0 \in \partial g(\vec{x})$, so we always have

$$- \nabla f(\overline{X})_i \in [-\lambda, \lambda] \quad \forall i \in I.$$

Suppose the following condition is true:

• (Strict complementarity) $\neg (\nabla f(\bar{X}))_i \in (-\lambda, \lambda)$ for all $i \in I$.

Then M is indeed an active C^p manifold for g at \overline{X} . Moreover, $(b \le)$ -regularity and strong (a)-regularity hold trivially for g along M at \overline{X} . Note that there is usually a bias between the center of the asymptotic distribution \overline{X} and the ground truth due to the regularization term.

2.4. Smoothly invertible maps and active manifolds. Performance of statistical estimation procedures strongly depends on the sensitivity of the problem to perturbation. A variety of estimation problems can in turn be modeled as the task of solving an inclusion $0 \in F(x)$ for some set-valued map F, whose values we can only approximate by sampling. We next review basic perturbation theory based on the inverse/implicit function theorem paradigm, while closely following the monograph [9]. A set-valued map $F : \mathbb{R}^d \Rightarrow \mathbb{R}^m$ is an assignment that maps a point $x \in \mathbb{R}^d$ to a set $F(x) \subseteq \mathbb{R}^d$

R^m. Set-valued maps always admit a set-valued inverse:

$$F^{-1}(y) = x : y \in F(x)$$
.

The domain and graph of *F* are defined, respectively, as

$$\operatorname{dom} F := x : F(x) = \emptyset$$
 and $\operatorname{gph} F := (x, y) : y \in F(x)$.

We will be interested in the sensitivity of the solutions to the system $v \in F(x)$ with respect to perturbations of the left-hand side V, or equivalently, the variational behavior of the map $V \to F^{-1}(V)$. In particular, we will be interested in settings when the graph of F^{-1} coincides locally around a base point (V, X) with a graph of a smooth map. This is the content of the following definition.

DEFINITION 2.6 (Smooth invertibility). Consider a set-valued map $F: \mathbb{R}^d \Rightarrow \mathbb{R}^m$ and a pair $(\overline{x}, \overline{v}) \in \operatorname{gph} F$. We say that F is C^p invertible around $(\overline{x}, \overline{v})$ with inverse $\sigma(\cdot)$ if there exists a single-valued C^p -smooth map $\sigma(\cdot)$ and a neighborhood U of $(\overline{v}, \overline{x})$ satisfying

$$U \cap gph F^{-1} = U \cap gph \sigma$$
.

The definition might seem odd at first: there is nothing "smooth" about F, and yet we require the graph of F^{-1} to coincide with a graph of a smooth function near $(\overline{v}, \overline{x})$. On the contrary, we will see that in a variety of settings this assumption is indeed valid. In particular, smooth invertibility is typical in nonlinear programming.

EXAMPLE 3 (Nonlinear programming). Returning to Example 1 with $p \ge 2$, define the set-valued map

$$F(x) = \nabla f(x) + N \times (x).$$

Then it is classically known that F is C^{p-1} invertible at $(\overline{X}, 0)$ if and only if the matrix

$$:= P \quad T_{M}(\bar{x}) \nabla_{xx}^{2} L(\bar{x}, \bar{y}) P_{T_{M}(\bar{x})}$$

is nonsingular on T_M (\bar{X}). In this case, the Jacobian of the inverse map is $\nabla \sigma$ (0) = † , where † denotes the Moore–Penrose pseudoinverse. It is worthwhile to note that can be equivalently written as P_{T_M} (\bar{X}) P_{T_M} (\bar{X}).

EXAMPLE 4 (1-regularization). Returning to Example 2 with $p \ge 2$, define the set-valued map $F(x) = \nabla f(x) + \lambda \partial (\cdot 1)(x)$. Then F is C^{p-1} invertible at $(\overline{X}, 0)$ if and only if the matrix $:= P \ T_M(\overline{X}) \nabla^2 f(\overline{X}) P_{T_M(\overline{X})}$ is nonsingular on $T_M(\overline{X})$.

Smooth invertibility is closely tied to active manifolds, and Example 3 and Example 4 are simple consequences. Indeed the following much more general statement is true. This result follows from a standard argument combining active manifolds and the implicit function theorem. The proof appears in Section 7.1 of the Supplementary Material [7]. We will require a mild assumption on the considered functions. Following [26], Definition 2.1, a function f is called $subdifferentially\ continuous$ at a point \overline{X} if for any sequences $(X_i, V_i) \in \operatorname{gph} \partial f$ converging to some pair $(\overline{X}, \overline{V}) \in \operatorname{gph} \partial f$, the function values $f(X_i)$ converge to $f(\overline{X})$. In particular, functions that are continuous on their domains and closed convex functions are subdifferentially continuous.

THEOREM 2.7 (Smooth invertibility and active manifolds). Consider the map

$$F(x) := A(x) + \partial f(x),$$

where $A: \mathbb{R}^d \to \mathbb{R}^d$ is C^p -smooth and $f: \mathbb{R}^d \to \mathbb{R} \cup \{\infty\}$ is subdifferentially continuous near a point \overline{X} . Suppose that f admits a C^{p+1} active manifold M at some point \overline{X} for $-A(\overline{X}) \in \partial f(\overline{X})$. Let G(X) = 0 be any C^{p+1} -smooth local defining equations for M near \overline{X} and let f be a C^{p+1} -smooth function that agrees with f on a neighborhood of \overline{X} in M. Define the map

$$H(x, y) := A(x) + \nabla \hat{f(x)} + \nabla G(x)$$

Then there exists a unique multiplier vector \overline{y} satisfying the condition $0 = H(\overline{x}, \overline{y})$. Moreover, F is C^p -invertible around $(0, \overline{x})$ with inverse σ (·) if and only if the matrix

$$:= P \quad T_{M}(\bar{x}) \nabla_{X} H(\bar{x}, \bar{y}) P_{T_{M}(\bar{x})}$$

is nonsingular on $T_M(\vec{x})$, and in this case equality $\nabla \sigma(0) = {}^{\dagger}$ holds.

3. Asymptotic normality of SAA. Before analyzing the asymptotic performance of stochastic approximation algorithms, it is instructive to first recall guarantees for sample average approximation (SAA), where the assumptions, conclusions, and arguments are much simpler to state. This is the content of this section: we derive the asymptotic distribution of the SAA estimator for nonsmooth problems. Throughout the section, we focus on the problem of finding a point *X* satisfying the variational inclusion

(3.1)
$$0 \in A(x) + H(x)$$
 where $A(x) = \mathop{\mathbb{E}}_{z \sim P} A(x, z)$.

Here $H: \mathbb{R}^d \Rightarrow \mathbb{R}^d$ is a set-valued map with closed graph, P is a fixed probability distribution on some measure space (Z, F), and $A: \mathbb{R}^d \times Z \to \mathbb{R}^d$ is a measurable map. We will impose the following assumption throughout the rest of the section.

Assumption A. The map F := A + H is C^1 -smoothly invertible near $(0, \overline{X})$ with inverse σ (·).

The SAA approach to solving (3.1) proceeds as follows. Let $S = \{z_1, \dots, R\}$ be i.i.d. samples drawn from P and let X_k be a solution of the problem

(3.2)
$$0 \in A s(x) + H(x)$$
 where $A s(x) := \frac{1}{k} \int_{i=1}^{k} A(x, z_i),$

assuming one exists. We will now show that the solutions of sample average approximations are asymptotically normal with covariance $\nabla \sigma(0) \cdot \text{Cov}(A(\bar{x}, z)) \cdot \nabla \sigma(0)$. Though variants of this result are known [13, 14, 29], we provide a short proof in Section 8.1 of the Supplementary Material [7] highlighting the use of the solution map $\sigma(\cdot)$. To this end, we impose the following standard assumption.

ASSUMPTION B (Integrability and smoothness). Suppose that there exists a neighborhood U around \overline{X} satisfying the following.

- 1. For almost every Z, the map $A(\cdot, z)$ is differentiable at every $X \in U$.
- 2. The second moment bounds hold:

$$\sup_{x \in U} \mathop{\mathsf{E}}_{z \sim P} A(x,\,z)^{-2} < \infty \qquad and \quad \sup_{x \in U} \mathop{\mathsf{E}}_{z \sim P} \nabla A(x,\,z)^{-2} \underset{\mathrm{op}}{\overset{2}{\sim}} < \infty.$$

The following theorem shows that as long as x_k eventually stay in a sufficiently small neighborhood of \overline{X} , the error $\overline{K}(X_k - \overline{X})$ is asymptotically normal with covariance $\nabla \sigma$ (0) $Cov(M(\bar{x}, z)) \cdot \nabla \sigma(0)$. Verifying that the problem (3.2) admits solutions x_k that are sufficiently close to \overline{X} is a separate and well-studied topic and we do not discuss it here.

THEOREM 3.1 (Sample average approximation). Suppose that Assumptions A and B hold. In particular, there exist $_{1}$, $_{2}$ > 0 and a $_{1}$ -smooth map σ : B $_{1}$ (0) \rightarrow B $_{2}$ ($_{2}$) with

$$gph \sigma = B_1(0) \times B_2(\overline{X}) \cap gph F^{-1}$$
.

Suppose moreover that there exists a square integrable function L(z) satisfying

(3.3)
$$\nabla A(x_1, z) - \nabla A(x_2, z) \le L(z)x_1 - x_2 \quad \forall x_1, x_2 \in B_2(\vec{x}).$$

Shrinking 2, if necessary, let us ensure that $2 \le \min\{\frac{\lim (\sigma)^{-1}}{2EL}, \frac{1}{2EL}\}$. Let $S = \{z_1, \dots, R\}$ be i.i.d. samples drawn from P and let X_k be a measurable selection of the solutions (3.2) such that $\Pr[X_k \in B_{-1}(\bar{X})] \to 1$ as K tends to infinity. Then the expansion holds: $\overline{K}(X_k - \bar{X}) = -\nabla \sigma (0) \cdot \overline{K} A(\bar{X}) - A S(\bar{X}) + OP(1),$

$$\sqrt[N]{k}(x_k - \overline{x}) = -\nabla\sigma(0) \cdot \sqrt[N]{k} A(\overline{x}) - A_S(\overline{x}) + o_P(1),$$

and therefore

(3.4)
$$\sqrt[4]{k(x_k - \overline{x})} \xrightarrow{D} N 0, \nabla \sigma (0) \cdot Cov A(\overline{x}, z) \cdot \nabla \sigma (0) .$$

Note that Theorem 3.1 assumes existence of a measurable selection of the solutions (3.2) such that $\Pr[x_k \in B_2(\bar{x})] \to 1$ as k tends to infinity. This is a very mild assumption and follows for example, from uniform convergence and smooth invertibility.

THEOREM 3.2 (Existence of measurable selections). Suppose that F is C^1 -smoothly invertible around near $(0, \overline{X})$ and that A(X, Z) and $\nabla A(X, Z)$ converge uniformly on some ball around \overline{X} , that is there exists 0 such that

$$\sup_{x \in B} \nabla A \, s(x) - \nabla A(x) = o_p(1) \quad and \quad \sup_{x \in B} \Delta s(x) - \Delta(x) = o_p(1).$$

Then there exists $\delta > 0$ and a measurable selection of the solutions (3.2) such that $\Pr[Xk \in B_{\delta}(\bar{X})] \to 1$ as K tends to infinity.

PROOF. Standard results on the implicit function theorem (see proof of [9], Theorem 3G.3) imply that there exist sufficiently small ε_2 , $\varepsilon_3 > 0$ such that whenever $\sup_{X \in B} (\bar{X}) A S(X) - A(X) < \varepsilon_2$ and $\sup_{X \in B} (\bar{X}) \nabla A S(X) - \nabla A(X) < \varepsilon_2$, the map $A_S + H$ is guaranteed to be smoothly invertible on $B_{\varepsilon_3}(\bar{X}) \times B \varepsilon_3(0)$. In particular, the solution $X_k \in B_{\varepsilon_3}(\bar{X})$ of (3.2) exists and is unique. To see measurability of X_k , observe that the maps A_S and $A_S + B_S$ and A_S

Our goal in the rest of the paper is to show that a simple online algorithm, namely the stochastic forward backward (SFB) method, under reasonable assumptions enjoys the same guarantees as (3.4) for SAA. Moreover, in the final section of the paper (Section 6), we will show that this performance is best possible among any estimation procedure in a local minimax sense, and therefore both SAA and SFB are asymptotically local minimax optimal.

4. Stochastic approximation: Assumptions & examples. We now move to stochastic approximation algorithms, and in this section set forth the algorithms we will consider and the relevant assumptions. The concrete examples we will present will all be geared to solving variational inclusions, but the specifics of this problem class is somewhat distracting. Therefore, we will instead only isolate the essential ingredients that are needed for our results to take hold. Setting the stage, our goal is to find a point ^X satisfying the inclusion

$$(4.1) 0 \in F(x),$$

where $F: \mathbb{R}^d \Rightarrow \mathbb{R}^d$ is a set-valued map. Throughout, we fix one such solution \overline{X} of (4.1). We will assume that in a certain sense, the problem (4.1) is "variationally smooth". That is, there exists a distinguished manifold M —the active manifold in concrete examples—containing \overline{X} and such that the map $X \to P$ $T_M(X)$ F(X) is single-valued and C^1 -smooth on M near \overline{X} . The following assumption makes this precise.

Assumption C (Smooth reduction). Suppose that there exists a C^p manifold $M \subset \mathbb{R}^d$ such that the following properties are true.

(C1) The map
$$F_M: M \to \mathbb{R}^d$$
 defined by
$$F_M(x) := P_{T_M(x)}F(x)$$

is single-valued on some neighborhood of \overline{X} in M .

(C2) There exists a neighborhood U of $(\overline{X}, 0)$ such that

$$U \cap gph F = U \cap gph(F_M + N_M).$$

We note that smooth invertibility of F can be easily characterized in terms of the covariant Jacobian $\nabla_M F_M(\vec{x})$. This is the content of the following lemma.

LEMMA 4.1 (Jacobian of the solution map). The map F is C^p -smoothly invertible around $(\overline{X}, 0)$ with localization $\sigma(\cdot)$ if and only if the linear map $P_{T_M}(\overline{X})\nabla F_M(\overline{X})P_{T_M}(\overline{X})$ is nonsingular on $T_M(\overline{X})$. In this case, the Jacobian of the localization is given by

$$\nabla \sigma \left(0 \right) = \ P_{T_M \left(\vec{X} \right)} \nabla_M \ F_M \left(\vec{X} \right) P_{T_M \left(\vec{X} \right)} \ ^{\dagger}.$$

PROOF. Let be a smooth extension of F to a neighborhood $V \subseteq \mathbb{R}^d$ of \overline{X} . In light of Assumption (C2), the graphs of F and +N M coincide near $(\overline{X}, 0)$, and therefore we can focus on existence of smooth localizations of $(+N M)^{-1}$. Applying Lemma 7.2 in the Supplementary Material [7] with $\overline{y} = 0$, we see that +N M is C^p -smoothly invertible around $(\overline{X}, 0)$ if and only if the linear map $P_{T_M(\overline{X})} \nabla (-\overline{X}) P_{T_M(\overline{X})}$ is nonsingular on $T_M(\overline{X})$. In this case, the Jacobian of the localization is given by $\nabla \sigma(0) = (P_{T_M(\overline{X})} \nabla (-\overline{X}) P_{T_M(\overline{X})})^{\dagger}$. Noting the equality $\nabla F_M(\overline{X}) P_{T_M(\overline{X})} = \nabla_M(-\overline{X}) P_{T_M(\overline{X})}$ completes the proof.

The stochastic approximation algorithms we consider assume access to a *generalized gradient mapping*:

G:
$$R_{++} \times R^d \times R^d \rightarrow R^d$$
.

Given $X_0 \in \mathbb{R}^d$, the algorithm iterates the update

(4.2)
$$X_{k+1} = X_k - \alpha_k G_{\alpha_k}(X_k, V_k),$$

where $\alpha_k > 0$ is a control sequence and V_k is stochastic noise. We will place relevant assumptions on the noise V_k later in Section 5.

We make two assumptions on G. The first (Assumption D) is similar to classical Lipschitz assumptions and ensures the steplength can only scale linearly in ν .

Assumption D (Steplength). We suppose that there exists a constant C > 0 and a neighborhood U of \bar{x} such that the estimate

$$\sup_{x \in U_{E}} G_{\alpha}(x, v) \leq C 1 + v ,$$

holds for all $v \in \mathbb{R}^d$ and $\alpha > 0$, where we set $U_F := U \cap \text{dom } F$.

The second assumption makes precise the relationship between the mapping $\,G\,$ and $\,F_M\,$.

ASSUMPTION E (Strong (a) and aiming). We suppose that there exist constants C, $\mu > 0$ and a neighborhood U of \overline{X} such that the following hold for all $\nu \in \mathbb{R}^d$ and $\alpha > 0$, where we set $U_F := U \cap \text{dom } F$.

(E1) (Tangent comparison) For all $x \in U_F$, we have

$$P_{T_M \; (P_M \; (x))} \; \; G_\alpha(x, \, \nu) - F \; \; P_M \; (x) \; \; - \nu \quad \leq C \; \; 1 + \nu \qquad ^2 \; \operatorname{dist}(x, \, M \;) + \alpha \; \; .$$

(E2) (Proximal Aiming) For $x \in U_F$, we have

$$G_{\alpha}(x, v) - v, x - P_{M}(x) \ge \mu \cdot \text{dist}(x, M) - 1 + v$$
 ² $o \cdot \text{dist}(x, M) + C\alpha$.

Some comments are in order. Assumption (E1) ensures that the direction of motion $G_{\alpha_k}(x_k, v_k)$ approximates well $F_M(P_M(x))$ in tangent directions to the manifold M. Assumption (E2) ensures that after subtracting the noise from $G_{\alpha_k}(x_k, v_k)$, the update direction $x_k - x_{k+1}$ locally points towards the manifold M. Note that the little-O term in (E2) depends only on dist(x, M) and not on O. We will later show that this ensures the iterates O approach the manifold O at a controlled rate.

4.1. Examples of stochastic approximation for variation inclusions. The rest of the section is devoted to examples of algorithms satisfying Assumptions D and E. In all cases, we will consider the task of solving the variational inclusion

$$(4.3) 0 \in A(x) + \partial g(x) + \partial f(x).$$

Here $A: \mathbb{R}^d \to \mathbb{R}^d$ is any single-valued continuous map, $f: \mathbb{R}^d \to \mathbb{R}^d$ is a closed function, and $g: \mathbb{R}^d \to \mathbb{R}^d$ is a closed function that is bounded from below. ² As explained in the Introduction, variational inclusions encompass a variety of problems, most-notably first-order optimality conditions for nonlinear programming and Nash equilibria of games. In order to identify (4.3) with (4.1), we define the set-valued map F to be

$$F(x) := A(x) + \partial g(x) + \partial f(x)$$
.

Throughout, we fix a point X satisfying the inclusion (4.3).

A classical algorithm for problem (4.3) is the stochastic forward–backward iteration, which proceeds by taking "forward-steps" on $A + \partial g$ and proximal steps on f . Specifically, given a current iterate X_t , the algorithm performs the update

(4.4) Choose
$$W_t \in \partial g(x_t)$$

$$Choose X_{t+1} \in \operatorname{prox}_{\alpha_t f} X_t - \alpha_t A(X_t) + W_t + V_t$$

where V_t is a noise sequence. The operator $G_\alpha(x, v)$ corresponding to this algorithm is simply

$$G_{\alpha}(x, v) := \frac{x - s_f(x - \alpha(A(x) + s_g(x) + v))}{\alpha},$$

where $S_g(X)$ is any selection of the subdifferential $\partial g(X)$ and $S_f(X)$ is any selection of the proximal map $\operatorname{prox}_{\alpha f}(X)$. The goal of this section is to verify Assumption E for this operator under a number of reasonable assumptions on A, g, and f.

In particular, the local boundedness condition D for G is widely used in the literature, with a variety of sufficient conditions known. The following lemma describes a number of such conditions, which we will use in what follows. The proof appear in Section 9.1 of the Supplementary Material [7].

LEMMA 4.2 (Local boundedness). Suppose that $A(\cdot)$ and $S_g(\cdot)$ are locally bounded around \overline{X} . Then Assumption D is guaranteed to hold in any of the following settings.

- 1. f is the indicator function of a closed set X.
- 2. f is convex and the function $X \to \text{dist}(0, \partial f(X))$ is bounded on dom f near \bar{X} .
- 3. f is Lipschitz continuous on dom $g \cap \text{dom } f$.

We next verify Assumption E in a number of reasonable settings; all proofs appear in the Supplementary Material [7]. In particular, it will be useful to note the following expression for F_M . We will use this lemma throughout the section, without explicit reference.

²In particular, $\operatorname{prox}_{\alpha f}(X)$ is nonempty for all $x \in \mathbb{R}^d$ and all $\alpha > 0$.

Lemma 4.3 (Local tangent reduction). Suppose that f and g are Lipschitz continuous on their domains, A is C^p -smooth, f+g admits an active C^{p+1} manifold at X for $\neg A(X)$, and f and g are both C^{p+1} -smooth and strongly (a) regular along M near X. Then Assumption G holds and G admits the simple form

$$(4.5) F_M(x) = P_{T_M(x)} A(x) + \nabla_M g(x) + \nabla_M f(x),$$

for all $x \in M$ near X.

4.1.1. *Stochastic forward algorithm* (f = 0). We begin with the simplest case of (4.3) where f = 0. In this case, the iteration (4.2) reduces to a pure stochastic forward algorithm and the map G takes the simple form

$$G_{\alpha}(x, v) := A(x) + s g(x) + v$$

which is independent of α . Let us introduce the following assumption on the problem data.

ASSUMPTION F (Assumptions for the forward algorithm). Suppose that f = 0 and that both $g(\cdot)$ and $A(\cdot)$ are Lipschitz continuous around \overline{X} . Suppose that $M \subseteq X$ is a C^p -smooth manifold for g at \overline{X} .

- (F1) (Strong (a)) The function g is strongly (a)-regular along M at \bar{X} .
- (F2) (Proximal aiming) There exists $\mu > 0$ such that the inequality holds:

$$(4.6) A(\overline{x}) + v, x - P_M(X) \ge \mu \cdot \operatorname{dist}(X, M) \text{for all } X \text{ near } \overline{X} \text{ and } v \in \partial g(X).$$

Note that Corollary 2.4 shows that the aiming condition (F2) holds as long as the inclusion $-A(\bar{x}) \in \partial g(\bar{x})$ holds, M is an active manifold for g at \bar{x} for $v = -A(\bar{x})$, and g is (b_{\leq}) -regular along M at \bar{x} . The following proposition shows that Assumption F suffices to ensure Assumption E. The proof appears in the Supplementary Material.

PROPOSITION 4.4 (Forward method). Assumption F implies Assumption E.

The following is now immediate.

COROLLARY 4.5 (Active manifolds). Suppose f = 0 and that both $g(\cdot)$ and $A(\cdot)$ are Lipschitz continuous around \overline{X} . Suppose moreover that the inclusion $\neg A(\overline{X}) \in \partial g(\overline{X})$ holds, that g admits a G^2 active manifold around \overline{X} for $\overline{V} = \neg A(\overline{X})$, and that g is both g -regular and strongly (a)-regular along g at g. Then Assumption g holds.

4.1.2. Stochastic projected forward algorithm ($f = \delta x$). Next, we focus on the particular instance of (4.3) where f is an indicator function of a closed set X. In this case, the iteration (4.2) reduces to a stochastic projected forward algorithm and the map G takes the form

$$G_{\alpha}(x,\,\nu):=\,\frac{x-sx\,(x-\alpha(A(x)+s\,\,g(x)+\nu))}{\alpha},$$

where $S_X(X)$ is a selection of the projection map $P_X(X)$. In order to ensure Assumption E for the stochastic projected forward method, we introduce the following assumption.

ASSUMPTION G (Assumptions for the projected gradient mapping). Suppose that f is the indicator function of a closed set X and both $g(\cdot)$ and $A(\cdot)$ are Lipschitz continuous around \overline{X} . Let $M \subseteq X$ be a C^2 manifold containing \overline{X} and suppose that f is C^2 on M near \overline{X} .

- (G1) (Strong (a)) The function g and set X are strongly (a)-regular along M at \bar{X} .
- (G2) (Proximal aiming) There exists $\mu > 0$ such that the inequality holds
- $(4.7) \qquad A(\overline{x}) + v, \ x P_M(X) \ge \mu \cdot \operatorname{dist}(X, M) \quad \forall x \in X \ near \ \overline{X} \ and \ v \in \partial g(x).$
 - (G3) (Condition (b)) The set X is (b_{\leq}) -regular along M at \overline{X} .

Note that a similar argument as Corollary 2.4 shows that the aiming condition (G2) holds as long as the inclusion $\neg A(\bar{x}) \in \hat{\partial}(g+f)(\bar{x})$ holds, M is an active manifold of g+f at \bar{x} for $v = \neg A(\bar{x})$, and g is (b_{\leq}) -regular along M at \bar{x} .

The following proposition shows that Assumption G is sufficient to ensure Assumption E.

PROPOSITION 4.6 (Projected forward method). *If Assumptions* D *and* G *hold, then so does Assumption* E.

The following is now immediate.

COROLLARY 4.7 (Active manifolds). Suppose that f is the indicator function of a closed set X and both $g(\cdot)$ and $A(\cdot)$ are Lipschitz continuous around \overline{X} . Suppose moreover the inclusion $-A(\overline{X}) \in \partial(g+f)(\overline{X})$ holds, g+f admits a C^2 active manifold around \overline{X} for the vector $\overline{V} = -A(\overline{X})$, and both g and g are g are g are g are and strongly g around g and g are g are g are g around g are g around g are g around g are g around g around g around g are g around g around

4.1.3. Stochastic forward–backward method (g = 0). Finally, we focus on the particular instance of (4.3) where g = 0. In this case, the iteration (4.2) reduces to a stochastic forward–backward algorithm and the map G becomes

$$G_{\alpha}(x, v) := \frac{x - s_f(x - \alpha(A(x) + v))}{\alpha}$$

In order to ensure Assumption E for the stochastic proximal gradient method, we introduce the following assumptions.

ASSUMPTION H (Assumptions for the forward–backward method). Suppose g=0 and $f(\cdot)$ and $A(\cdot)$ are Lipschitz continuous on dom f near \overline{X} . Suppose moreover that there exists a C^2 manifold $M\subset X$ containing \overline{X} and such that f is C^2 -smooth on M near \overline{X} .

- (H1) (Strong (a)) The function f is strongly (a)-regular along M at \bar{X} .
- (H2) (Proximal Aiming) There exists $\mu > 0$ such that the inequality

(4.8)
$$A(\bar{x}) + v, x - P_M(x) \ge \mu \cdot \text{dist}(x, M) - 1 + v \quad \text{o dist}(x, M)$$

holds for all $x \in \text{dom} f$ near \overline{X} and $v \in \partial f(x)$.

Note that Corollary 2.4 shows that the aiming condition (H2) holds as long as the inclusion $-A(\overline{x}) \in \partial f(\overline{x})$ holds, M is an active manifold for f at \overline{x} for $v = -A(\overline{x})$, and f is (b_{\leq}) -regular along M at \overline{x} .

The following proposition shows that Assumption H is sufficient to ensure Assumption E.

PROPOSITION 4.8 (Forward–backward method). *If Assumptions* D *and* H *hold, then so does Assumption* E.

The following is now immediate.

COROLLARY 4.9 (Active manifolds). Suppose g = 0 and both f and $A(\cdot)$ are Lipschitz continuous on dom^f near \overline{X} . Suppose moreover the inclusion— $A(\overline{X}) \in \hat{\partial} f(\overline{X})$ holds. Suppose that f admits a C^2 active manifold around \overline{X} for $\overline{V} = -A(\overline{X})$ and f is both $(b)_{\leq}$ -regular and strongly (a)-regular along M at \overline{X} . Then Assumption E holds.

5. Asymptotic normality. Next, we impose two assumptions on the step-size α_k and the noise sequence α_k . The first is standard, and is summarized next.

ASSUMPTION I (Standing assumptions). Assume the following.

- (I1) The map G is measurable.
- (I2) There exist constants c_1 , $c_2 > 0$ and $y \in (1/2, 1]$ such that

$$\frac{c_1}{k^{\gamma}} \le \alpha_k \le \frac{c_2}{k^{\gamma}}.$$

(I3) $\{Vk\}$ is a martingale difference sequence w.r.t. to the increasing sequence of σ -fields

$$F_k = \sigma(x_j : j \le k \text{ and } v_j : j < k),$$

and there exists a function $q: \mathbb{R}^d \to \mathbb{R}_+$ that is bounded on bounded sets with

$$E[v_k|F_k] = 0$$
 and $E v_k^4|F_k < q(x_k)$.

We let $E_k[\cdot] = E[\cdot|F_k]$ denote the conditional expectation.

(I4) The inclusion $X_k \in \text{dom } F$ holds for all $k \ge 1$.

All items in Assumption I are standard in the literature on stochastic approximation methods and mirror for example those found in [8], Assumption C. The only exception is the fourth moment bound on v_k , which stipulates that v_k has slightly lighter tails. This bound appears to be necessary for the setting we consider.

To prove our asymptotic normality results, we impose a further assumption on the noise sequence V_k , which also appears in [12], Assumption D. Before stating it, as motivation, consider the stochastic variational inequality (4.3) given by:

$$0 \in A(x) + \partial f(x) + \partial g(x)$$
 where $A(x) = \mathop{\mathbb{E}}_{z \sim P} A(x, z)$.

Then the noise V_k in the algorithm (4.4) takes the form

$$V_k = A(x_k; z_k) - A(x_k).$$

Equivalently, we may decompose the right-hand side as

$$v_k = A(\overline{x}; z_k) - A(\overline{x}) + A(x_k; z_k) - A(x_k; z_k) + A(x_k) + A(x_k) - A(x_k)$$

$$=: v_k^{(1)}$$

The two components $V_k^{(1)}$ and $V_k^{(2)}(X_k)$ are qualitatively different in the following sense. On one hand, the sum $\frac{1}{\sqrt{K}} \int_{i=1}^{k} V_i^{(1)}$ clearly converges to a zero-mean normal vector as long as the covariance $\text{Cov}(A(\overline{X}, Z))$ exists. On the other hand, $V_k^{(2)}(X_k)$ is small in the sense that $\text{E}_k V_k^{(2)}(X_k)^2 \leq 2 \cdot \text{E}_z[L(z)^2] \cdot X_k - \overline{X}^2$, where L(Z) is a Lipschitz constant of $A(\cdot, Z)$. With this example in mind, we introduce the following assumption on the noise sequence.

ASSUMPTION J. Fix a point $\bar{x} \in \text{dom } F$ at which Assumption C holds and let U be a matrix whose column vectors form an orthogonal basis of $[M](\bar{x})$. Recall that $E_k[\cdot]$ denote the conditional expectation with respect to F_k . We suppose the noise sequence has decomposable structure $V_k = V_k^{(1)} + V_k^{(2)}(X_k)$, where $V_k^{(2)}$: $\text{dom } F \to R$ d is a random function satisfying

$$E_k \ U \ V_k^{(2)}(X)^2 \le C_X - X^2 \ \text{for all } X \in \text{dom } F \ \text{near } X,$$

and some C > 0. In addition, we suppose that for all $x \in \text{dom } F$, we have $E_k[v_k^{(1)}] = E_k[v_k^{(2)}(x)] = 0$ and the following limit holds:

$$\frac{1}{\sqrt{K}} \int_{i=1}^{K} U v_i^{(1)} \xrightarrow{D} N 0, U U$$

for some symmetric positive semidefinite matrix

Note that Assumption J only requires that $V_k^{(1)}$ and $V_k^{(2)}$ have zero conditional mean, which is weaker than being independent of the previous iterates. We are now ready to state the main result of this work—asymptotic normality for stochastic approximation algorithms.

THEOREM 5.1 (Asymptotic normality). Suppose that Assumption C, D, E, I, and J hold. Suppose that $y \in (\frac{1}{2}, 1)$ and that the sequence x_k generated by the process (4.2) converges to \overline{x} with probability one. Suppose that there exists a constant $\mu > 0$ satisfying

(5.1)
$$\nabla_M F_M(\vec{x}) v, v \ge \mu v^{-2} \quad \text{for all } v \in T_M(\vec{x}).$$

Then F is C^p -smoothly invertible around $(\overline{X}, 0)$ with inverse $\sigma(\cdot)$ and the average iterate $\overline{X}_k = \frac{1}{k} \int_{i=1}^k X_i$ admits the expansion

$$\sqrt[4]{k}(\bar{X_k} - \bar{X}) = -\sqrt[4]{\frac{1}{k}} \int_{i=1}^{k} U U \nabla_M F_M(\bar{X}) U^{-1} U V_i^{(1)} + o_P(1),$$

and hence

$$\sqrt[4]{\overline{k}(\overline{X_k}-\overline{X})} \stackrel{D}{\longrightarrow} N \ 0, \ \nabla\sigma \ (0) \cdot \cdot \nabla\sigma \ (0) \ .$$

Moreover, $\nabla \sigma$ (0) can be equivalently written as $\nabla \sigma$ (0) = $(P_{T_M}(\vec{x})\nabla_M F_M(\vec{x})P_{T_M}(\vec{x}))^{\dagger}$.

The conclusion of this theorem is surprising: although the sequence X_k never reaches the manifold, the limiting distribution of $\overline{k}(\overline{X_k}-\overline{X})$ is supported on the tangent space $T_M(\overline{X})$. Thus asymptotically, the "directions of nonsmoothness," which are normal to M, are quickly "averaged out." When $G_{\alpha_k}(X_k, V_k)$ is bounded away from 0 for all K, this means that X_k must oscillate across the manifold, instead of approaching it from one direction.

5.1. Asymptotic normality in nonlinear programming. As a simple illustration of Theorem 5.1, we now spell out the consequence for the stochastic projected gradient method for stochastic nonlinear programming, already discussed in Example 1. Namely, consider the problem (2.7) and let \overline{X} be a local minimizer. Suppose that G are G-smooth near \overline{X} and G takes the form G and each function G is G-smooth near G. Consider the following stochastic projected gradient method:

Sample:
$$Z_k \sim P$$
,

(5.2) Update:
$$X_{k+1} \in Px \mid X_k - \alpha_k \nabla f(x \mid k; z_k)$$
.

In order to understand the asymptotics of the algorithm, as in Example 1, let \overline{Y} be the Lagrange multiplier vector and suppose that LICQ and strict complementarity holds. Suppose moreover the second-order sufficient conditions: there exists $\mu > 0$ such that

$$(5.3) W \nabla_{xx}^2 L(\overline{x}, \overline{y}) w \ge \mu w ^2 \text{for all } w \in T_M(\overline{x}).$$

Note that, as explained in Example 3, this condition is simply the requirement that the covariant Hessian of $f := f_0 + \delta x$

$$\nabla_{M}^{2} f(\overline{x}) = P_{T_{M}(\overline{x})} \nabla_{xx}^{2} L x , y P_{T_{M}(\overline{x})}$$

is positive definite on $T_M(\vec{x})$. Finally, to ensure our noise sequence

$$\begin{aligned} v_k &= \nabla f\left(x \mid k; z_k\right) - \nabla f\left(x \mid k\right) \\ &= \nabla f\left(\overline{x}; z_k\right) - \nabla f\left(\overline{x}\right) + \underbrace{\nabla f\left(x \mid k; z_k\right) - \nabla f\left(\overline{x}; z_k\right) + \nabla f\left(\overline{x}\right) - \nabla f\left(x \mid k\right)}_{=:v_k^{(2)}}, \\ &= v_k^{(1)} \end{aligned}$$

satisfies Assumptions I and J, we assume the stochasticity is sufficiently well behaved:

ASSUMPTION K (Stochastic gradients). As a function of X, the fourth moment

$$x \in X \to \mathbb{E}_{z \sim P} \quad \nabla f(x; z) - \nabla f(x)$$

is bounded on bounded sets. Moreover, there exists C > 0 such that

$$\mathsf{E}_{z\sim P} \ \nabla f(x;z) - \nabla f(\bar{x};z)^2 \leq Cx - \bar{x}^2 \text{ for all } x\in X.$$

Finally, the gradients $P_{T_M(\vec{x})} \nabla f(\vec{x}; z)$ have finite covariance = $Cov(P_{T_M(\vec{x})} \nabla f(\vec{x}; z))$.

With these assumptions in hand, we have the following asymptotic normality result for nonlinear programming—a direct corollary of Theorem 5.1.

COROLLARY 5.2 (Asymptotic normality in nonlinear programming). Suppose that LICQ, strict complementary, second-order sufficient conditions, and Assumption K hold. Suppose that $y \in (\frac{1}{2}, 1)$ and consider the iterates x_k generated by the stochastic projected gradient method (5.2). Then if x_k converges to \overline{x} with probability 1, the average iterate $\overline{X_k} = \frac{1}{K} \quad \underset{i=1}{\overset{k}{\bigvee}} X_i$ admits the expansion

$$\sqrt[4]{k}(\overline{x_{k}} - \overline{x}) = -\sqrt[4]{\frac{1}{k}} \int_{i=1}^{k} U U \nabla_{xx}^{2} L(\overline{x_{i}}, \overline{y}) U^{-1} U v_{i}^{(1)} + o_{P}(1),$$

where the columns of U form an orthonormal basis of $T_M(\bar{X})$. Consequently, asymptotic normality holds:

$$\sqrt[4]{\overline{k}}(\overline{X_k}-\overline{X}) \stackrel{d}{\longrightarrow} N \ 0, \ \nabla\sigma \ (0) \cdot \text{Cov} \ \nabla f \ (\overline{X}; z) \cdot \nabla\sigma \ (0) \quad ,$$

where $\nabla \sigma (0) = (P_{T_M(\vec{x})} \nabla^2_{xx} L(\vec{x}, \vec{y}) P_{T_M(\vec{x})})^{\dagger}$.

As stated in the Introduction, this appears to be the first asymptotic normality guarantee for the standard stochastic projected gradient method in general nonlinear programming problems with C^3 data, even in the convex setting. Finally, we note that even for simple optimization problems, dual averaging procedures can achieve suboptimal convergence [12]. This is surprising since such methods identify the active manifold [17] (also see [12], Section 4.1), while projected stochastic gradient methods do not.

It is instructive to look at three problem formulations for sparse recovery: EXAMPLE 5.

(regularized)
$$\min_{x \in \mathcal{X}} E f(x, z) + \lambda x$$
₁,

$$(I_1 \text{ constraint})$$
 $\min_{\substack{x \in A}} E f(x, z)$,

$$(I_{1} \text{ constraint}) \qquad \min_{\substack{x \\ 1 \le A}} \ \mathsf{E} \ f(x, z) \ ,$$

$$(I_{0} \text{ constraint}) \qquad \min_{\substack{|\sup(x)| \le s}} \ \mathsf{E} \ f(x, z) \ .$$

Problem (regularized) is typically solved by the stochastic proximal algorithm, while $(I_1 \text{ constraint})$ and $(I_0 \text{ constraint})$ are solved by the stochastic projected gradient method. Both methods are trivially examples of the algorithm (4.4) that we have studied in the section. Let us now look at the asymptotic covariance of these methods corresponding to the three problems. To this end, let X denote the optimal solution for the three problems and suppose that X y = A and $|\sup(X)| = S$. Without loss of generality suppose moreover $\sup(X) = \{1, \ldots, S\}$ In all cases, under the regularity conditions discussed in the section, the asymptotic covariance of the average iterate is

$$\nabla \sigma (0) \cdot \operatorname{Cov} \nabla f \times_{\mathcal{I}} Z \cdot \nabla \sigma (0)$$
.

Thus the only distinction is in Jacobian of the solution map $\nabla \sigma$ (0). It is straightforward to see that the active manifold (under strict complementarity) for (regularized) and (l_0 constraint) is

$$M_{1,3} = \mathbb{R}^{s} \times \{0\}^{d-s}$$

while the active manifold for $(l_1 \text{ constraint})$ is

$$M_2 = M_{1,3} \cap x : |x_i| = A$$
.

Because the active manifold M in all cases is piecewise linear, an application of Theorem 2.7 yields the expression:

$$\nabla \sigma (0) = P_{T_M(x)} \to \nabla^2 f(x), z P_{T_M(x)}^{\dagger}.$$

For the problem (regularized) and (I_0 constraint), the tangent space is simply $T_{M_{1.3}}(X) = \mathbb{R}^s \times \{0\}^{d-s}$, while for (I_1 constraint), the tangent space is smaller

$$T_{M_2} X = v \in T_{M_{1,3}} X : \underset{i=1}{\text{sign } X_i} v_i = 0$$
.

In particular, the asymptotic covariance corresponding to (I_1 constraint) is no larger in the Loewner order than that of (regularized) and (I_2 constraint). Consequently, the formulation (I_2 constraint) may be preferable when I_3 as I_4 is known. The caveat, however, is that LASSO solution is biased due to the regularization term and therefore a comparison of the three formulation purely based on the asymptotic covariance is not entirely justified.

6. Asymptotic optimality of SAA and SFB. In this section, we show that the asymptotic covariance in (3.4) is the best possible among all estimators of \overline{X} , and therefore both SAA and SFB are asymptotically optimal. Namely, we will lower-bound the performance of *any* estimation procedure for finding a solution of an adversarially-chosen sequence of small perturbations of the target problem. In order to specify this sequence, define the set

$$G:= g: Z \to \mathbb{R}^d: \underset{z \sim P}{\mathsf{E}} g(z) = 0, \underset{z \sim P}{\mathsf{E}} g(z)^2 < \infty$$

Fix now a function $g \in G$ and an arbitrary C^3 -smooth function $h: \mathbb{R} \to [-1, 1]$ such that its first three derivatives are bounded and h(t) = t for all $t \in [-1/2, 1/2]$. Now for each $u \in \mathbb{R}^d$, define a new probability distribution D^u whose density is given by

(6.1)
$$dP_{u}(z) := \frac{1 + h(u - g(z))}{C(u)} dP(z),$$

where C(u) is the normalizing constant $C(u) := 1 + h(u \ g(z)) \ dP(z)$. Thus each vector $u \in \mathbb{R}^d$ induces the perturbed problem

(6.2)
$$0 \in L(x, u) + H(x)$$
 where $L(x, u) = \mathop{\mathbb{E}}_{z \sim P_u} A(x, z)$.

Reassuringly, the following lemma shows that map(x, u) $\to L(x, u)$ is C^1 near (\overline{x} , 0). All proofs of results in this section appear in Section 11 of the Supplementary Material [7].

LEMMA 6.1. The map
$$(x, u) \to L(x, u)$$
 is C^1 near $(\overline{X}, 0)$ with partial derivatives $\nabla_X L(\overline{X}, 0) = \nabla A(\overline{X})$ and $\nabla_u L(X, 0) = \mathop{\mathbb{E}}_{z \sim P} A(\overline{X}, z) g(z)$.

The family of problems (6.2) would not be particularly useful if their solution would vary wildly in U. On the contrary, the following lemma shows that for all small U, each problem (6.2) admits a unique solution in U, which moreover varies smoothly in U. We will use the following standard notation. A map $\sigma(\cdot)$ is called a *localization* of a set-valued map $F(\cdot)$ around a pair $(\overline{u}, \overline{v}) \in \operatorname{gph} F$ if the two sets, $\operatorname{gph} \sigma$ and $\operatorname{gph} F$, coincide locally around $(\overline{u}, \overline{v})$.

LEMMA 6.2 (Derivative of the solution map). *The solution map*

$$S(u) = x : 0 \in L(x, u) + H(x)$$

admits a single-valued localization $S(\cdot)$ around around $(0, \overline{X})$ that is differentiable at 0 with Jacobian

$$\nabla s(0) = -\nabla \sigma(0) \cdot \mathop{\mathbb{E}}_{z \sim P} A(\overline{X}, z) g(z) \quad .$$

In light of Lemma 6.2, for all small u, we define the solution $\overline{X}_u := s(u)$. The following theorem provides an asymptotic lower bound on the performance of any estimator when applied to the problems within our parametric family. We let $E_{P_u^k}$ denote the expectation with respect to k i.i.d. observations $z_i \sim P_u$.

THEOREM 6.3 (Local minimax). Let $L: \mathbb{R}^d \to [0, \infty)$ be symmetric, quasiconvex, and lower semicontinuous, let $X_k: Z^k \to U$ be a sequence of estimators, and set $g(z) := A(\overline{X}, z) - A(\overline{X})$. Then the inequality

(6.3)
$$\sup_{I \subset \mathbb{R}^d, |I| < \infty} \liminf_{k \to \infty} \max_{u \in I} \mathsf{E}_{P_{u'}^k \setminus \overline{k}} L \sqrt{\overline{k}} (x_k - \overline{x}_{u'}^{\vee} \setminus \overline{k}) \ge \mathsf{E} L(Z)$$

holds, where $Z \sim N(0, \nabla \sigma(0) \cdot \text{Cov}(A(\bar{x}, z)) \cdot \nabla \sigma(0)$).

In particular, applying Theorem 6.3 with quadratics L yields a lower bound on the achievable covariance among any estimator. We will now show that both SAA and SFB fulfill (6.3) with equality, and therefore in a precise sense *asymptotically minimax optimal*. Note that we already know that the asymptotic covariance $\sigma(0) \cdot \text{Cov}(A(\overline{x}, z)) \cdot \nabla \sigma(0)$ is achieved by both SAA (Theorem 3.1) and SFB (Theorem 5.1) when applied to the *fixed problem* u = 0. It remains therefore to argue that $K(x_k - \overline{x}_u)$ along the perturbed sequence of problems is asymptotically independent of U. This is the content of the following theorem.

THEOREM 6.4 (Tightness of SAA). Under the same assumptions as Theorem 3.1, the sample average approximation estimator $x_k := x_k$ satisfies (6.3) with equality for any bounded continuous function $L: \mathbb{R}^d \to [0, \infty)$.

SFB enjoys completely analogous results, which we summarize next.

THEOREM 6.5 (Tightness of SFB). Suppose the same setting as Theorem 5.1 and that $V_i^{(1)} = A(\overline{X}, Z_i) - A(\overline{X})$ with $Z_i \sim \text{iid } P$ and such that $E_i = V_i^{(1)} =$

THEOREM 6.6 (Tightness of SFB for nonlinear programming). *Under the same assumptions as Corollary* 5.2, *the average iterate* $^{X}k := \frac{1}{k} \int_{i=1}^{k} ^{X_i} satisfies$ (6.3) *with equality for any bounded continuous function* $L : \mathbb{R}^d \to [0, \infty)$.

We note that asymptotic optimality of SFB for smooth problems was proved in [4], Theorem 5.6, and the proof we present of the three theorems above is an adaptation of the argument therein.

Acknowledgments. The authors thank John Duchi for insightful comments about the paper and encouraging us to prove that SFB attains the asymptotic lower bound in (6.3).

Funding. Research of D. Davis supported by an Alfred P. Sloan research fellowship and NSF Grant DMS-2047637.

Research of D. Drusvyatskiy was supported by NSF Grants DMS-1651851, NSF DMS-2306322, and CCF-2023166.

SUPPLEMENTARY MATERIAL

Supplement to "Asymptotic normality and optimality in nonsmooth stochastic approximation" (DOI: 10.1214/24-AOS2401SUPP; .pdf). Supplementary information.

REFERENCES

- [1] BORWEIN, J. and L EWIS, A. S. (2017). Convex Analysis and Nonlinear Optimization: Theory and Examples. Springer, Berlin.
- [2] BOUMAL, N. (2023). An Introduction to Optimization on Smooth Manifolds. Cambridge Univ. Press, Cambridge. MR4533407
- [3] CLARKE, F. H., LEDYAEV, Y. S., STERN, R. J. and WOLENSKI, P. R. (1998). Nonsmooth Analysis and Control Theory. Graduate Texts in Mathematics 178. Springer, New York. MR1488695
- [4] CUTLER, J., DÍAZ, M. and DRUSVYATSKIY, D. (2024). Stochastic approximation with decision-dependent distributions: Asymptotic normality and optimality. *J. Mach. Learn. Res.* 25 Paper No. 90, 49. MR4749126
- [5] DAVIS, D., DRUSVYATSKIY, D. and CHARISOPOULOS, V. (2019). Stochastic algorithms with geometric step decay converge linearly on sharp functions. arXiv preprint. Available at arXiv:1907.09547.
- [6] DAVIS, D., DRUSVYATSKIY, D. and JIANG, L. (2021). Active manifolds, stratifications, and convergence to local minima in nonsmooth optimization. arXiv preprint. Available at arXiv:2108.11832v2.
- [7] DAVIS, D., DRUSVYATSKIY, D. and JIANG, L. (2024). Supplement to "Asymptotic normality and optimality in nonsmooth stochastic approximation." https://doi.org/10.1214/24-AOS2401SUPP
- [8] DAVIS, D., DRUSVYATSKIY, D., KAKADE, S. and LEE, J. D. (2020). Stochastic subgradient method converges on tame functions. *Found. Comput. Math.* 20 119–154. MR4056927 https://doi.org/10.1007/s10208-018-09409-5
- [9] DONTCHEV, A. L. and ROCKAFELLAR, R. T. (2009). Implicit Functions and Solution Mappings: A view from variational analysis. Springer Monographs in Mathematics. Springer, Dordrecht. MR2515104 https://doi.org/10.1007/978-0-387-87821-8
- [10] Drusvyatskiy, D., Ioffe, A. D. and Lewis, A. S. (2016). Generic minimizing behavior in semialgebraic optimization. *SIAM J. Optim.* **26** 513–534. MR3461323 https://doi.org/10.1137/15M1020770
- [11] DRUSVYATSKIY, D. and L EWIS, A. S. (2014). Optimality, identifiability, and sensitivity. *Math. Program.* **147** 467–498. MR3258532 https://doi.org/10.1007/s10107-013-0730-4
- [12] DUCHI, J. C. and R UAN, F. (2021). Asymptotic optimality in stochastic optimization. Ann. Statist. 49 21–48. MR4206668 https://doi.org/10.1214/19-AOS1831
- [13] Dupačová, J. and Wets, R. (1988). Asymptotic behavior of statistical estimators and of optimal solutions of stochastic optimization problems. *Ann. Statist.* **16** 1517–1549. MR0964937 https://doi.org/10.1214/aos/1176351052
- [14] KING, A. J. and ROCKAFELLAR, R. T. (1993). Asymptotic theory for solutions in statistical estimation and stochastic programming. *Math. Oper. Res.* **18** 148–162. MR1250111 https://doi.org/10.1287/moor.18. 1.148

- [15] LE CAM, L. and YANG, G. L. (2000). Asymptotics in Statistics: Some Basic Concepts, 2nd ed. Springer Series in Statistics. Springer, New York. MR1784901 https://doi.org/10.1007/978-1-4612-1166-2
- [16] LEE, J. M. (2013). Smooth manifolds. In *Introduction to Smooth Manifolds*, 2nd ed. *Graduate Texts in Mathematics* 218 1–31. Springer, New York. MR2954043
- [17] LEE, S. and WRIGHT, S. J. (2012). Manifold identification in dual averaging for regularized stochastic online learning. J. Mach. Learn. Res. 13 1705–1744. MR2956341
- [18] LEMARÉCHAL, C., OUSTRY, F. and SAGASTIZÁBAL, C. (2000). The *U*-Lagrangian of a convex function. *Trans. Amer. Math. Soc.* **352** 711–729. MR1487623 https://doi.org/10.1090/S0002-9947-99-02243-6
- [19] LEWIS, A. S. (2002). Active sets, nonsmoothness, and sensitivity. SIAM J. Optim. 13 702–725. MR1972212 https://doi.org/10.1137/S1052623401387623
- [20] LEWIS, A. S. and M ALICK, J. (2008). Alternating projections on manifolds. *Math. Oper. Res.* 33 216–234. MR2393548 https://doi.org/10.1287/moor.1070.0291
- [21] MIFFLIN, R. and SAGASTIZÁBAL, C. (2005). A *VU*-algorithm for convex minimization. *Math. Program.* **104** 583–608. MR2179252 https://doi.org/10.1007/s10107-005-0630-3
- [22] MILLER, S. A. and MALICK, J. (2005). Newton methods for nonsmooth convex minimization: Connections among *U*-Lagrangian, Riemannian Newton and SQP methods. *Math. Program.* 104 609–633. MR2179253 https://doi.org/10.1007/s10107-005-0631-2
- [23] MORDUKHOVICH, B. S. (2006). Variational Analysis and Generalized Differentiation. I: Basic Theory. Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences] 330. Springer, Berlin. MR2191744
- [24] NESTEROV, Y. (2009). Primal-dual subgradient methods for convex problems. *Math. Program.* 120 221–259. MR2496434 https://doi.org/10.1007/s10107-007-0149-x
- [25] PENOT, J.-P. (2012). Calculus Without Derivatives. Graduate Texts in Mathematics 266. Springer, New York. MR2986672 https://doi.org/10.1007/978-1-4614-4538-8
- [26] POLIQUIN, R. A. and R OCKAFELLAR, R. T. (1996). Prox-regular functions in variational analysis. *Trans. Amer. Math. Soc.* **348** 1805–1838. MR1333397 https://doi.org/10.1090/S0002-9947-96-01544-9
- [27] POLYAK, B. T. and JUDITSKY, A. B. (1992). Acceleration of stochastic approximation by averaging. SIAM J. Control Optim. 30 838–855. MR1167814 https://doi.org/10.1137/0330046
- [28] ROCKAFELLAR, R. T. and WETS, R. J.-B. (2009). Variational Analysis 317. Springer, Berlin.
- [29] Shapiro, A. (1989). Asymptotic properties of statistical estimators in stochastic programming. Ann. Statist. 17 841–858. MR0994271 https://doi.org/10.1214/aos/1176347146
- [30] SHAPIRO , A. (2003). On a class of nonsmooth composite functions. *Math. Oper. Res.* 28 677–692. MR2015908 https://doi.org/10.1287/moor.28.4.677.20512
- [31] VAN DER VAART, A. W. (2000). Asymptotic Statistics. Cambridge Series in Statistical and Probabilistic Mathematics 3. Cambridge Univ. Press, Cambridge. MR1652247 https://doi.org/10.1017/CBO9780511802256
- [32] WRIGHT, S. J. (1993). Identifiable surfaces in constrained optimization. SIAM J. Control Optim. 31 1063– 1079. MR1227547 https://doi.org/10.1137/0331048
- [33] XIAO, L. (2009). Dual averaging method for regularized stochastic learning and online optimization. In *Advances in Neural Information Processing Systems* **22**.