Convergence of a Stochastic Subgradient Method with Averaging for Nonsmooth Nonconvex Constrained Optimization*

Andrzej Ruszczyński[†]

December 16, 2019

Abstract

We prove convergence of a single time-scale stochastic subgradient method with subgradient averaging for constrained problems with a nonsmooth and nonconvex objective function having the property of generalized differentiability. As a tool of our analysis, we also prove a chain rule on a path for such functions.

Keywords: Stochastic Subgradient Method, Nonsmooth Optimization, Generalized Differentiable Functions, Chain Rule

1 Introduction

We consider the problem

$$\min_{x \in V} f(x) \tag{1}$$

where $X \subset \mathbb{R}^n$ is convex and closed, and $f : \mathbb{R}^n \to \mathbb{R}$ is a Lipschitz continuous function, which may be neither convex nor smooth. The subgradients of $f(\cdot)$ are not available; instead, we postulate access to their random estimates.

Research on stochastic subgradient methods for nonsmooth and nonconvex functions started in late 1970's. Early contributions are due to Nurminski, who considered weakly convex functions and established a general methodology for studying convergence of non-monotonic methods [20], Gupal and his co-authors, who considered convolution smoothing (mollification) of Lipschitz functions and resulting finite-difference methods [11], and Norkin, who considered unconstrained problems with "generalized differentiable" functions [17, Ch. 3 and 7]. Recently, by an approach via differential inclusions, Duchi and Ruan [10] studied proximal methods for sum-composite problems with weakly convex functions, Davis *et al.* [8] proved convergence of the subgradient method for locally Lipschitz Whitney stratifiable functions with constraints, and Majewski *et al.* [15] studied several methods for subdifferentially regular Lipschitz functions.

Our objective is to show that a single time-scale stochastic subgradient method with direction averaging [21, 22], is convergent for a broad class of functions enjoying the property of "generalized differentiability," which contains all classes of functions mentioned above, as well as their compositions.

^{*}This publication was supported by the NSF Award DMS-1312016.

[†]Rutgers University, Department of Management Science and Information Systems, Piscataway, NJ 08854, USA; email: rusz@rutgers.edu

Our analysis follows the approach of relating a stochastic approximation algorithm to a continuous-time dynamical system, pioneered in [14, 13] and developed in many works (see, *e.g.*, [12] and the references therein). Extension to multifunctions was proposed in [1] and further developed, among others, in [3, 10, 8, 15].

For the purpose of our analysis, we also prove a chain rule on a path under generalized differentiability, which may be of independent interest.

We illustrate the use of the method for training a *ReLu* neural network.

2 The chain formula on a path

Norkin [19] introduced the following class of functions.

Definition 2.1. A function $f : \mathbb{R}^n \to \mathbb{R}$ is differentiable in a generalized sense at a point $x \in \mathbb{R}^n$, if an open set $U \subset \mathbb{R}^n$ containing x, and a nonempty, convex, compact valued, and upper semicontinuous multifunction $G_f : U \rightrightarrows \mathbb{R}^n$ exist, such that for all $y \in U$ and all $g \in G_f(y)$ the following equation is true:

$$f(y) = f(x) + \langle g(y), y - x \rangle + o(x, y, g),$$

with

$$\lim_{y \to x} \sup_{g \in G(y)} \frac{o(x, y, g)}{\|y - x\|} = 0.$$

The set $G_f(y)$ is the generalized subdifferential of f at y. If a function is differentiable in a generalized sense at every $x \in \mathbb{R}^n$ with the same generalized subdifferential mapping $G_f : \mathbb{R}^n \rightrightarrows \mathbb{R}^n$, we call it differentiable in a generalized sense.

A function $f : \mathbb{R}^n \to \mathbb{R}^m$ is differentiable in a generalized sense, if each of its component functions, $f_i : \mathbb{R}^n \to \mathbb{R}, i = 1, ..., m$, has this property.

The class of such functions is contained in the set of locally Lipschitz functions [17, Thm. 1.1], and contains all subdifferentially regular functions [5], Whitney stratifiable Lipschitz functions [9], semismooth functions [16], and their compositions. In fact, if a function is differentiable in generalized sense and has directional derivatives at x in every direction, then it is semismooth at x. The Clarke subdifferential $\partial f(x)$ is an inclusion-minimal generalized subdifferential, but the generalized subdifferential mapping $G_f(\cdot)$ is not uniquely defined in Definition 2.1, which plays a role in our considerations. For stochastic optimization, essential is the closure of the class of such functions with respect to expectation, which allows for easy generation of stochastic subgradients. In the Appendix we recall basic properties of functions differentiable in a generalized sense. For thorough exposition, see [17, Ch. 1 and 6].

Our interest is in a formula for calculating the increment of a function $f : \mathbb{R}^n \to \mathbb{R}$ along a path $p : [0,\infty) \to \mathbb{R}^n$, which is at the core of the analysis of nonsmooth and stochastic optimization algorithms (see [9, 7] and the references therein). For an absolutely continuous function $p : [0,\infty) \to \mathbb{R}^n$ we denote by $\dot{p}(\cdot)$ its weak derivative, that is, a measurable function such that

$$p(t) = p(0) + \int_0^t \dot{p}(s) \, ds, \quad \forall t \ge 0.$$

Theorem 2.2. If $f : \mathbb{R}^n \to \mathbb{R}$ and $p : [0, \infty) \to \mathbb{R}^n$ are differentiable in a generalized sense, then for every T > 0, any generalized subdifferential $G_f(\cdot)$, and every selection $g(p(t)) \in G_f(p(t))$, we have

$$f(p(T)) - f(p(0)) = \int_0^T \left\langle g(p(t)), \dot{p}(t) \right\rangle dt.$$
 (2)

Proof. Consider the function

$$\varphi(\varepsilon) = \int_0^T f(p(t+\varepsilon)) dt, \quad \varepsilon \ge 0.$$

Its right derivative at 0 can be calculated in two ways:

$$\varphi_{+}'(0) = \lim_{\varepsilon \downarrow 0} \frac{1}{\varepsilon} \left[\varphi(\varepsilon) - \varphi(0) \right] = \lim_{\varepsilon \downarrow 0} \frac{1}{\varepsilon} \left[\int_{0}^{T} f(p(t+\varepsilon)) dt - \int_{0}^{T} f(p(t)) dt \right]$$

$$= \lim_{\varepsilon \downarrow 0} \frac{1}{\varepsilon} \left[\int_{\varepsilon}^{T+\varepsilon} f(p(\tau)) d\tau - \int_{0}^{T} f(p(t)) dt \right]$$

$$= \lim_{\varepsilon \downarrow 0} \frac{1}{\varepsilon} \left[\int_{T}^{T+\varepsilon} f(p(t)) dt - \int_{0}^{\varepsilon} f(p(t)) dt \right] = f(p(T)) - f(p(0)).$$

$$(3)$$

On the other hand,

$$\varphi'_{+}(0) = \lim_{\varepsilon \downarrow 0} \int_{0}^{T} \frac{1}{\varepsilon} \left[f(p(t+\varepsilon)) - f(p(t)) \right] dt.$$
(4)

By the generalized differentiability of $f(\cdot)$, the differential quotient under the integral can be expanded as follows:

$$\frac{1}{\varepsilon} \Big[f(p(t+\varepsilon)) dt - f(p(t)) \Big]
= \frac{1}{\varepsilon} \Big\langle g(p(t+\varepsilon)), p(t+\varepsilon) - p(t) \Big\rangle + \frac{1}{\varepsilon} o\big(p(t), p(t+\varepsilon), g(p(t+\varepsilon)) \big),$$
(5)
with
$$\lim_{\varepsilon \downarrow 0} \frac{1}{\varepsilon} o\big(p(t), p(t+\varepsilon), g(p(t+\varepsilon)) \big) = 0.$$

Since $p(\cdot)$ is differentiable in a generalized sense, it is locally Lipschitz continuous [17, Thm. 1.1], hence absolutely continuous. Thus, for almost all *t*, we have $\frac{1}{\varepsilon} [p(t+\varepsilon) - p(t)] = \dot{p}(t) + r(t,\varepsilon)$, with $\lim_{\varepsilon \downarrow 0} r(t,\varepsilon) = 0$. Combining it with (5), and using the local boundedness of generalized gradients, we obtain

$$\frac{1}{\varepsilon} \left[f(p(t+\varepsilon)) - f(p(t)) \right] = \left\langle g(p(t+\varepsilon)), \dot{p}(t) \right\rangle + O(t,\varepsilon), \tag{6}$$

with $\lim_{\epsilon \downarrow 0} O(t, \epsilon) = 0$. By [17, Thm. 1.6] (Theorem A.1), the function $\psi(t) = f(p(t))$ is differentiable in a generalized sense as well and

$$G_{\boldsymbol{\Psi}}(t) = \left\{ \langle g, h \rangle : g \in G_f(p(t)), \ h \in G_p(t) \right\}$$

is its generalized subdifferential. By virtue of [17, Cor. 1.5] (Theorem A.3), any generalized subdifferential mapping $G_{\psi}(\cdot)$ is single-valued except for a countable number of points in [0,1]. Since it is upper semicontinuous, it is continuous almost everywhere. By [17, Thm. 1.12] (Theorem A.2), almost everywhere $G_p(t) = \{\dot{p}(t)\}$. Then for any $h(t + \varepsilon) \in G_p(t + \varepsilon)$ and for almost all t,

$$\lim_{\varepsilon \downarrow 0} \left\langle g(p(t+\varepsilon)), h(t+\varepsilon) \right\rangle = \left\langle g(p(t)), \dot{p}(t) \right\rangle.$$

Therefore, for almost all t,

$$\lim_{\varepsilon \downarrow 0} \left\langle g(p(t+\varepsilon)), \dot{p}(t) \right\rangle = \left\langle g(p(t)), \dot{p}(t) \right\rangle + \lim_{\varepsilon \downarrow 0} \left\langle g(p(t+\varepsilon)), \dot{p}(t) - h(t+\varepsilon) \right\rangle = \left\langle g(p(t)), \dot{p}(t) \right\rangle,$$

where the last equation follows from the local boundedness of $G_f(\cdot)$ and the continuity of $G_p(\cdot)$ at the points of differentability. Thus, for almost all *t*, we can pass to the limit in (6):

$$\lim_{\varepsilon \downarrow 0} \frac{1}{\varepsilon} \left[f(p(t+\varepsilon)) \, dt - f(p(t)) \right] = \left\langle g(p(t)), \dot{p}(t) \right\rangle.$$

We can now use the Lebesgue theorem and pass to the limit under the integral in (4):

$$\varphi'_{+}(0) = \int_{0}^{T} \lim_{\varepsilon \downarrow 0} \frac{1}{\varepsilon} \left[f(p(t+\varepsilon)) \, dt - f(p(t)) \right] \, dt = \int_{0}^{T} \left\langle g(p(t)), \dot{p}(t) \right\rangle \, dt.$$

Comparison with (3) yields (2).

3 The single time-scale method with subgradient averaging

We briefly recall from [21, 22] a stochastic approximation algorithm for solving problem (1) where only random estimates of subgradients of f are available.

The method generates two random sequences: approximate solutions $\{x^k\}$ and path-averaged stochastic subgradients $\{z^k\}$, defined on a certain probability space (Ω, \mathscr{F}, P) . We let \mathscr{F}_k to be the σ -algebra generated by $\{x^0, \ldots, x^k, z^0, \ldots, z^k\}$. We assume that for each k, we can observe an \mathscr{F}_k -measurable random vector $g^k \in \mathbb{R}^n$, such that, for some \mathscr{F}_k -measurable vector r^k , we have $g^k - r^k \in G_f(x^k)$. Further assumptions on the errors r^k will be specified in section 4.

The method proceeds for k = 0, 1, 2... as follows (a > 0 and $\beta > 0$ are fixed parameters). We compute

$$y^{k} = \underset{y \in X}{\operatorname{argmin}} \left\{ \langle z^{k}, y - x^{k} \rangle + \frac{\beta}{2} \|y - x^{k}\|^{2} \right\},$$

$$(7)$$

and, with an \mathscr{F}_k -measurable stepsize $\tau_k \in (0, \min(1, 1/a)]$, we set

$$x^{k+1} = x^k + \tau_k (y^k - x^k).$$
(8)

Then we observe g^{k+1} at x^{k+1} , and update the averaged stochastic subgradient as

$$z^{k+1} = (1 - a\tau_k)z^k + a\tau_k g^{k+1}.$$
(9)

Convergence of the method was proved in [22] for weakly convex functions $f(\cdot)$. Unfortunately, this class does not contain functions with downward cusps, which are common in modern machine learning models (see section 5).

4 Convergence analysis

We call a point $x^* \in \mathbb{R}^n$ *Clarke stationary* of problem (1), if

$$0 \in \partial f(x^*) + N_X(x^*), \tag{10}$$

where $N_X(x^*)$ denotes the normal cone to X at x^* . The set of Clarke stationary points of problem (1) is denoted by X^* .

We start from a useful property of the gap function $\eta: X \times \mathbb{R}^n \to (-\infty, 0]$,

$$\eta(x,z) = \min_{y \in X} \left\{ \langle z, y - x \rangle + \frac{\beta}{2} \| y - x \|^2 \right\}.$$
(11)

We denote the minimizer in (11) by $\bar{y}(x,z)$. Since it is a projection of $x - z/\beta$ on X, we observe that

$$\langle z, \bar{y}(x,z) - x \rangle + \beta \| \bar{y}(x,z) - x \|^2 \le 0.$$
⁽¹²⁾

Moreover, a point $x^* \in X^*$ if and only if $g^* \in \partial f(x^*)$ exists such that $\eta(x^*, g^*) = 0$.

We analyze convergence of the algorithm (7)–(9) under the following conditions, the first three of which are assumed to hold with probability 1:

- (A1) All iterates x^k belong to a compact set;
- (A2) $\tau_k \in (0, \min(1, 1/a)]$ for all k, $\lim_{k\to\infty} \tau_k = 0$, $\sum_{k=0}^{\infty} \tau_k = \infty$;
- (A3) For all $k, r^k = e^k + \delta^k$, with $\sum_{k=0}^{\infty} \tau_k e^k$ convergent, and $\lim_{k\to\infty} \delta_k = 0$;

(A4) The set $\{f(x) : x \in X^*\}$ does not contain an interval of nonzero length.

Condition (A3) can be satisfied for a martingale $\sum_{k=0}^{\infty} \tau_k e^k$, but can also hold for broad classes of dependent "noise" sequences $\{e^k\}$ [12]. Condition (A4) is true for Whitney stratifiable functions [2, Cor. 5], but we need to assume it here.

We have the following elementary property of the sequence $\{z^k\}$.

Lemma 4.1. Suppose the sequence $\{x^k\}$ is included in a set $A \subset \mathbb{R}^n$ and conditions (A2) and (A3) are satisfied. Then

$$\lim_{k \to \infty} \operatorname{dist}(z^k, B) = 0, \quad where \quad B = \operatorname{conv}\left(\bigcup_{x \in A} \partial f(x)\right)$$

Proof. Using (A2), we define the quantities $\tilde{z}^k = z^k + a \sum_{i=k}^{\infty} \tau_i e^i$ and establish the recursive relation

$$ilde{z}^{k+1} = (1-a au_k) ilde{z}^k + a au_k g^k + au_k \Delta_k, \quad k=0,1,2,\ldots,$$

where $g^k \in B$ and $\Delta_k = a\delta^k + a\sum_{j=k}^{\infty} \tau_j e^j \to 0$ a.s.. The convexity of the distance function and (A2) yield the result.

Theorem 4.2. If assumptions (A1)–(A4) are satisfied, then, with probability 1, every accumulation point \hat{x} of the sequence $\{x^k\}$ is Clarke stationary, and the sequence $\{f(x^k)\}$ is convergent.

Proof. Due to (A1), by virtue of Lemma 4.1, the sequence $\{z^k\}$ is bounded. We divide the proof into three standard steps.

Step 1: The Limiting Dynamical System. We define $p^k = (x^k, z^k)$, accumulated stepsizes $t_k = \sum_{j=0}^{k-1} \tau_j$, k = 0, 1, 2, ..., and we construct the interpolated trajectory

$$P_0(t) = p^k + \frac{t - t_k}{\tau_k} (p^{k+1} - p^k), \quad t_k \le t \le t_{k+1}, \quad k = 0, 1, 2, \dots$$

For an increasing sequence of positive numbers $\{s_k\}$ diverging to ∞ , we define shifted trajectories $P_k(t) = P_0(t + s_k)$. Recall that $P_k(t) = (X_k(t), Z_k(t))$.

By [15, Thm. 3.2], for any increasing sequence $\{n_k\}$ of positive integers, there exist a subsequence $\{\tilde{n}_k\}$ and absolutely continuous functions $X_{\infty} : [0, +\infty) \to X$ and $Z_{\infty} : [0, +\infty) \to \mathbb{R}^n$ such that for any T > 0

$$\lim_{k\to\infty}\sup_{t\in[0,T]}\left(\left\|X_{\tilde{n}_k}(t)-X_{\infty}(t)\right\|+\left\|Z_{\tilde{n}_k}(t)-Z_{\infty}(t)\right\|\right)=0,$$

and $(X_{\infty}(\cdot), Z_{\infty}(\cdot))$ is a solution of the system of differential equations and inclusions:

$$\dot{x}(t) = \bar{y}(x(t), z(t)) - x(t),$$
(13)

$$\dot{z}(t) \in a\big(\partial f(x(t)) - z(t)\big). \tag{14}$$

Moreover, for any $t \ge 0$, the pair $(X_{\infty}(t), Z_{\infty}(t))$ is an accumulation point of the sequence $\{(x^k, z^k)\}$.

Step 2: Descent Along a Path. We use the Lyapunov function

$$W(x,z) = af(x) - \eta(x,z)$$

For any solution (X(t),Z(t)) of the system (13)–(14), and for any T > 0, we estimate the difference W(X(T),Z(T)) - W(X(0),Z(0)). We split $W(X(\cdot),Z(\cdot))$ into a generalized differentiable composition $f(X(\cdot))$ and the "classical" part $\eta(X(\cdot),Z(\cdot))$.

Since the path $X(\cdot)$ satisfies (13) and $\bar{y}(\cdot, \cdot)$ is continuous, $X(\cdot)$ is continuously differentiable. Thus, we can use Theorem 2.2 to conclude that for any $g(X(\cdot)) \in \partial f(X(\cdot))$

$$f(X(T)) - f(X(0)) = \int_0^T \left\langle g(X(t)), \dot{X}(t) \right\rangle dt = \int_0^T \left\langle g(X(t)), \bar{y}(X(t), Z(t)) - X(t) \right\rangle dt.$$
(15)

On the other hand, since $\bar{y}(x,z)$ is unique, the function $\eta(\cdot, \cdot)$ is continuously differentiable. Therefore, the chain formula holds for it as well:

$$\eta(X(T),Z(T)) - \eta(X(0),Z(0)) = \int_0^T \left\langle \nabla_x \eta(X(t),Z(t)), \dot{X}(t) \right\rangle dt + \int_0^T \left\langle \nabla_z \eta(X(t),Z(t)), \dot{Z}(t) \right\rangle dt$$

Substituting $\nabla_x \eta(x,z) = -z + \beta(x - \bar{y}(x,z)), \nabla_z \eta(x,z) = \bar{y}(x,z) - x$ and $\dot{Z}(t) = a(\hat{g}(X(t)) - Z(t))$ with some $\hat{g}(X(\cdot)) \in \partial f(X(\cdot))$, and using (12) we obtain

$$\begin{split} &\eta(X(T), Z(T)) - \eta(X(0), Z(0)) \\ &= \int_0^T \left\langle -Z(t) + \beta(X(t) - \bar{y}(X(t), Z(t))), \, \bar{y}(X(t), Z(t)) - X(t) \right\rangle dt \\ &\quad + a \int_0^T \left\langle \bar{y}(X(t), Z(t)) - X(t), \, \hat{g}(X(t)) - Z(t) \right\rangle dt \\ &\geq a \int_0^T \left\langle \bar{y}(X(t), Z(t)) - X(t), \, \hat{g}(X(t)) - Z(t) \right\rangle dt \\ &\geq a \int_0^T \left\langle \bar{y}(X(t), Z(t)) - X(t), \, \hat{g}(X(t)) \right\rangle dt + a\beta \int_0^T \left\| \bar{y}(X(t), Z(t)) - X(t) \right\|^2 dt. \end{split}$$

We substitute the subgradient selector $g(X(t)) = \hat{g}(X(t))$ into (15) and combine it with the last inequality, concluding that

$$W(X(T), Z(T)) - W(X(0), Z(0)) \le -a\beta \int_0^T \left\| \bar{y}(X(t), Z(t)) - X(t) \right\|^2 dt = -a\beta \int_0^T \left\| \dot{X}(t) \right\|^2 dt.$$
(16)

Step 3: Analysis of Limit Points. Define the set $\mathscr{S} = \{(x,z) \in X^* \times \mathbb{R}^n : \eta(x,z) = 0\}$. Suppose (\bar{x},\bar{z}) is an accumulation point of the sequence $\{(x^k, z^k)\}$. If $\eta(\bar{x}, \bar{z}) < 0$, then every solution (X(t), Z(t)) of the system (13)–(14), starting from $(X(0), Z(0)) = (\bar{x}, \bar{z})$ has $\dot{X}(0) \neq 0$. Using (16) and arguing as in [10, Thm. 3.20] or [15, Thm. 3.5], we obtain a contradiction. Therefore, we must have $\eta(\bar{x}, \bar{z}) = 0$. Suppose $\bar{x} \notin X^*$. Then

$$\operatorname{dist}\left(0,\partial f(\bar{x}) + N_X(\bar{x})\right) > 0. \tag{17}$$

Suppose $X(t) = \bar{x}$ for all $t \ge 0$. The inclusion (14) simplifies: $\dot{z}(t) \in a(\partial f(\bar{x}) - z(t))$. By using the convex Lyapunov function $V(z) = \text{dist}(z, \partial f(\bar{x}))$ and applying the classical chain formula on the path $Z(\cdot)$ [4], we deduce that

$$\lim_{t \to \infty} \operatorname{dist} \left(Z(t), \partial f(\bar{x}) \right) = 0.$$
(18)

It follows from (17)–(18) that T > 0 exists, such that $-Z(T) \notin N_X(\bar{x})$, which yields $\dot{X}(T) \neq 0$. Consequently, the path X(t) starting from \bar{x} cannot be constant. But then again T > 0 exists, such that $\dot{X}(T) \neq 0$. By Step 1, the pair (X(T), Z(T)) would have to be an accumulation point of of the sequence $\{(x^k, z^k)\}$, a case already excluded. We conclude that every accumulation point (\bar{x}, \bar{z}) of the sequence $\{(x^k, z^k)\}$ is in \mathscr{S} . The convergence of the sequence $\{W(x^k, z^k)\}$ then follows in the same way as [10, Thm. 3.20] or [15, Thm. 3.5]. As $\eta(x^k, z^k) \to 0$, the convergence of $\{f(x^k)\}$ follows as well.

Directly from Lemma 4.1 we obtain convergence of averaged stochastic subgradients.

Corollary 4.3. If the sequence $\{x^k\}$ is convergent to a single point \bar{x} , then every accumulation point of $\{z^k\}$ is an element of $\partial f(\bar{x})$.

5 Example

A *Rectified Linear Unit (ReLU)* neural network [18] predicts a random outcome $Y \in \mathbb{R}^m$ from random features $X \in \mathbb{R}^n$ by a nonconvex nonsmooth function y(X, W), defined recursively as follows:

$$s_1 = X$$
, $s_{\ell+1} = (W_\ell s_\ell)_+$, $\ell = 1, 2, \dots, L-1$, $y(X, W) = W_L s_L$

where $(v)_+ = \max(0, v)$, componentwise. The decision variables are $W_1, \ldots, W_{L-1} \in \mathbb{R}^{n \times n}$ and $W_L \in \mathbb{R}^{m \times n}$. The simplest training problem is:

$$\min_{W \in \mathscr{W}} f(W) \stackrel{\triangle}{=} \frac{1}{2} \mathbb{E} \big[\| y(X, W) - Y \|^2 \big], \tag{19}$$

where \mathscr{W} is a box about 0. The function f(W) is not subdifferentially regular. It is not Whitney stratifiable, in general, because this property is not preserved under the expected value operator. However, we can use Theorems A.1 and A.4 to verify that it is differentiable in a generalized sense, and to calculate its stochastic subgradients. For a random data point (X^k, Y^k) we subdifferentiate the function under the expected value in (19) by recursive application of Theorem A.1. In particular, for L = 2 and m = 1 we have y(X, W) = $W_2(W_1X)_+$, and $g^k = (y(X^k, W^k) - Y^k) [D^k(W_2^k)^T (X^k)^T (W_1^k X^k)_+^T]$. Here, D^k is a diagonal matrix with 1 on position (i, i), if $(W_1^k X^k)_i > 0$, and 0 otherwise. A typical run of the stochastic subgradient method and the method with direction averaging is shown in Fig. 1, on an example of predicting wine quality [6], with identical random starting points, sequences of observations, and schedules of stepsizes: $\tau_k =$ 0.03/(1 + 5k/N), where N = 500,000. The coefficient a = 0.1. For comparison, the loss of a simple regression model is 666.



Figure 1: Comparison of methods with (lower graph) and without averaging (upper graph).

References

- [1] M. Benaïm, J. Hofbauer, and S. Sorin. Stochastic approximations and differential inclusions. *SIAM Journal on Control and Optimization*, 44(1):328–348, 2005.
- [2] J. Bolte, A. Daniilidis, A. Lewis, and M. Shiota. Clarke subgradients of stratifiable functions. SIAM Journal on Optimization, 18(2):556–572, 2007.
- [3] V. S. Borkar. Stochastic Approximation: a Dynamical Systems Viewpoint. Springer, New York, 2009.
- [4] H. Brézis. Monotonicity methods in Hilbert spaces and some applications to nonlinear partial differential equations. In *Contributions to Nonlinear Functional Analysis*, pages 101–156. Elsevier, 1971.
- [5] F. H. Clarke. Generalized gradients and applications. *Transactions of the American Mathematical Society*, 205:247–262, 1975.
- [6] P. Cortez, A. Cerdeira, F. Almeida, T. Matos, and J. Reis. Modeling wine preferences by data mining from physicochemical properties. *Decision Support Systems*, 47(4):547–553, 2009.
- [7] D. Davis and D. Drusvyatskiy. Stochastic model-based minimization of weakly convex functions. *SIAM Journal on Optimization*, 29(1):207–239, 2019.
- [8] D. Davis, D. Drusvyatskiy, S. Kakade, and J. D. Lee. Stochastic subgradient method converges on tame functions. *Foundations of Computational Mathematics*, pages 1–36, 2018.
- [9] D. Drusvyatskiy, A. D. Ioffe, and A. S. Lewis. Curves of descent. SIAM Journal on Control and Optimization, 53(1):114–138, 2015.

- [10] J. C. Duchi and F. Ruan. Stochastic methods for composite and weakly convex optimization problems. SIAM Journal on Optimization, 28(4):3229–3259, 2018.
- [11] A. M. Gupal. Stochastic Methods for Solving Nonsmooth Extremal Problems. Naukova Dumka, Kiev, 1979.
- [12] H. Kushner and G. G. Yin. Stochastic Approximation Algorithms and Applications. Springer, New York, 2003.
- [13] H. J. Kushner and D. S. Clark. Stochastic Approximation Methods for Constrained and Cnconstrained Systems. Springer, New York, 1978.
- [14] L. Ljung. Analysis of recursive stochastic algorithms. *IEEE Transactions on Automatic Control*, 22(4):551–575, 1977.
- [15] S. Majewski, B. Miasojedow, and E. Moulines. Analysis of nonsmooth stochastic approximation: the differential inclusion approach. arXiv preprint arXiv:1805.01916, 2018.
- [16] R. Mifflin. Semismooth and semiconvex functions in constrained optimization. SIAM Journal on Control and Optimization, 15(6):959–972, 1977.
- [17] V. S. Mikhalevich, A. M. Gupal, and V. I. Norkin. *Nonconvex Optimization Methods*. Nauka, Moscow, 1987.
- [18] V. Nair and G. E. Hinton. Rectified linear units improve restricted boltzmann machines. In Proceedings of the 27th International Conference on Machine Learning (ICML-10), pages 807–814, 2010.
- [19] V. I. Norkin. Generalized-differentiable functions. *Cybernetics and Systems Analysis*, 16(1):10–12, 1980.
- [20] E. A. Nurminski. *Numerical Methods for Solving Deterministic and Stochastic Minimax Problems*. Naukova Dumka, Kiev, 1979.
- [21] A. Ruszczyński. A method of feasible directions for solving nonsmooth stochastic programming problems. In F. Archetti, G. Di Pillo, and M. Lucertini, editors, *Stochastic Programming*, pages 258–271. Springer, 1986.
- [22] A. Ruszczyński. A linearization method for nonsmooth stochastic programming problems. Mathematics of Operations Research, 12(1):32–49, 1987.

Appendix A Generalized differentiability of functions

Compositions of generalized diifferentiable functions are crucial in our analysis.

Theorem A.1. [17, Thm. 1.6] If $h : \mathbb{R}^m \to \mathbb{R}$ and $f_i : \mathbb{R}^n \to \mathbb{R}$, i = 1, ..., m, are differentiable in a generalized sense, then the composition $\Psi(x) = h(f_1(x), ..., f_m(x))$ is differentiable in a generalized sense, and at any point $x \in \mathbb{R}^n$ we can define the generalized subdifferential of Ψ as follows:

$$G_{\psi}(x) = \operatorname{conv}\left\{g \in \mathbb{R}^{n} : g = \begin{bmatrix}g_{1} & \cdots & g_{m}\end{bmatrix}g_{0}, \\ with \ g_{0} \in G_{h}\left(f_{1}(x), \dots, f_{m}(x)\right) \text{ and } g_{j} \in G_{f_{j}}(x), \ j = 1, \dots, m\right\}.$$
(20)

Even if we take $G_h(\cdot) = \partial h(\cdot)$ and $G_{f_j}(\cdot) = \partial f_j(\cdot)$, j = 1, ..., m, we may obtain $G_{\psi}(\cdot) \neq \partial \psi(\cdot)$, but G_{ψ} defined above satisfies Definition 2.1.

Theorem A.2. [17, Thm. 1.12] If $f : \mathbb{R}^n \to \mathbb{R}$ is differentiable in a generalized sense, then for almost all $x \in \mathbb{R}^n$ we have $G_f(x) = \{\nabla f(x)\}$.

Functions of one variable have the following remarkable property.

Theorem A.3. [17, Cor. 1.5] If $f : \mathbb{R} \to \mathbb{R}$ is differentiable in a generalized sense, then the set of points x at which a generalized subdifferential $G_f(x)$ is not a singleton is at most countable.

For stochastic optimization, essential is the closure of the class functions differentiable in a generalized sense with respect to expectation.

Theorem A.4. [17, Thm. 23.1] Suppose (Ω, \mathscr{F}, P) is a probability space and a function $f : \mathbb{R}^n \times \Omega \to \mathbb{R}$ is differentiable in a generalized sense with respect to x for all $\omega \in \Omega$, and integrable with respect to ω for all $x \in \mathbb{R}^n$. Let $G_f : \mathbb{R}^n \times \Omega \Rightarrow \mathbb{R}^n$ be a multifunction, which is measurable with respect to ω for all $x \in \mathbb{R}^n$, and which is a generalized subdifferential mapping of $f(\cdot, \omega)$ for all $\omega \in \Omega$. If for every compact set $K \subset \mathbb{R}^n$ an integrable function $L_K : \Omega \to \mathbb{R}$ exists, such that $\sup_{x \in K} \sup_{g \in G_f(x, \omega)} ||g|| \le L_K(\omega), \omega \in \Omega$, then the function

$$F(x) = \int_{\Omega} f(x, \omega) P(d\omega), \quad x \in \mathbb{R}^n,$$

is differentiable in a generalized sense, and the multifunction

$$G_F(x) = \int_{\Omega} G_f(x, \omega) P(d\omega), \quad x \in \mathbb{R}^n,$$

is its generalized subdifferential mapping.