

Failure of uniform laws of large numbers for subdifferentials and beyond

Lai Tian*

Johannes O. Royset*

November 20, 2025

Abstract

We provide counterexamples showing that uniform laws of large numbers do not hold for subdifferentials under natural assumptions. Our results apply to random Lipschitz functions and random convex functions with a finite number of smooth pieces. Consequently, they resolve the questions posed by Shapiro and Xu [J. Math. Anal. Appl., 325(2), 2007] in the negative and highlight the obstacles nonsmoothness poses to uniform results.

1 Introduction

Uniform laws of large numbers (LLNs) for gradients are of fundamental importance in stochastic programming [20, Section 7.2.5], and their validity and rate of convergence are now well understood for random smooth functions; see, e.g., [9, 5]. However, the question of whether similar uniform laws hold for their nonsmooth counterparts has remained open for nearly two decades [19, Remark 2]. Positive results are known only for special function classes [17, 19, 21, 24, 15], for *enlarged* subdifferentials [18, 19], and for weaker notions of convergence [12, 4, 16]; see also the discussion in [20, Section 7.2.6] and [10, p. 471].

To set the stage, let nonempty $X \subset \mathbb{R}^d$ be compact, $\Xi \subset \mathbb{R}^m$, and ξ^1, ξ^2, \dots be independent and identically distributed (iid) Ξ -valued random variables on the complete probability space $(\Omega, \mathcal{F}, \mathbb{P})$. (Throughout, we use boldface for random variables.) Given $f : \Xi \times \mathbb{R}^d \rightarrow \mathbb{R}$ with measurable $f(\cdot, x)$ for all $x \in \mathbb{R}^d$, assume the expectation function $x \mapsto \mathbb{E}[f(\xi, x)]$ is finite valued and $f(\xi, \cdot)$ is $L(\xi)$ -Lipschitz on X for any $\xi \in \Xi$, with measurable L and $\mathbb{E}[L(\xi)] < \infty$. A general uniform LLN for *enlarged* set-valued mappings appears in [19, Theorem 2] and can be applied to the Clarke subdifferential $\partial_x f(\xi, x) = \partial(f(\xi, \cdot))(x)$. Specifically, for any fixed $r > 0$, \mathbb{P} -almost surely, one has

$$\lim_{\nu \rightarrow \infty} \sup_{x \in X} \text{exs} \left(\frac{1}{\nu} \sum_{i=1}^{\nu} \partial_x f(\xi^i, x); \bigcup_{y \in \mathbb{B}(x, r) \cap X} \mathbb{E}[\partial_x f(\xi, y)] \right) = 0, \quad (1)$$

$$\lim_{\nu \rightarrow \infty} \sup_{x \in X} \text{exs} \left(\mathbb{E}[\partial_x f(\xi, x)]; \bigcup_{y \in \mathbb{B}(x, r) \cap X} \frac{1}{\nu} \sum_{i=1}^{\nu} \partial_x f(\xi^i, y) \right) = 0, \quad (2)$$

where $\text{exs}(\cdot; \cdot)$ is the *excess* of a set relative to another one and the integral is taken in the Aumann sense; see the end of this section for a summary of notation. When $f(\xi, \cdot)$ is further assumed to be

*Daniel J. Epstein Department of Industrial and Systems Engineering, University of Southern California, Los Angeles, CA. Emails: {laitian, royset}@usc.edu

subdifferentially regular [14, Definition 7.25] on X for all $\xi \in \Xi$, we have

$$\mathbb{E}[\partial_x f(\boldsymbol{\xi}, x)] = \partial \mathbb{E}[f(\boldsymbol{\xi}, \cdot)](x), \quad \frac{1}{\nu} \sum_{i=1}^{\nu} \partial_x f(\xi^i, x) = \partial \left(\frac{1}{\nu} \sum_{i=1}^{\nu} f(\xi^i, \cdot) \right)(x),$$

for any $\xi^1, \dots, \xi^\nu \in \Xi$ and $x \in X$; see [3, Theorem 2.7.2]. It is then natural to ask, as also explicitly mentioned in [19, Remark 2], whether we can set $r = 0$ in (1) and/or (2), which, if both are true, would imply the following uniform law for subdifferentials in the Hausdorff sense:

$$\lim_{\nu \rightarrow \infty} \sup_{x \in X} d \left(\mathbb{E}[\partial_x f(\boldsymbol{\xi}, x)], \frac{1}{\nu} \sum_{i=1}^{\nu} \partial_x f(\xi^i, x) \right) = 0. \quad (3)$$

While the general questions about random Lipschitz functions remain elusive, the uniform law in (3) is actually achievable for several special function classes. One sufficient condition is that, for every $x \in X$, the subdifferential $\partial_x f(\xi, x)$ is a singleton for almost every $\xi \in \Xi$, which implies (3) and also that $x \mapsto \mathbb{E}[f(\boldsymbol{\xi}, x)]$ is continuously differentiable; see [17, Proposition 2.2], [19, Theorem 6], and [20, Theorem 7.52]. Beyond requiring the expectation function to be smooth, the work [21] guarantees (3) when the subdifferential mapping $\partial_x f$ satisfies a separable range condition; see [21, Theorem 4] and also [10, Theorem 5.1.31]. While this may be appealing for certain applications, the range of $\partial_x f$ is not separable even for the simple random convex function $(\xi, x) \mapsto f(\xi, x) = \max\{x - \xi, 0\}$ with $X = [0, 1]$ and ξ^i uniformly distributed on $\Xi = [0, 1]$ as pointed out in [12, Example 3.4]; see also [10, Example 5.1.33]. The recent work [15] establishes the uniform law (3) for a subclass of random convex-composite functions. The uniform law in [15, Theorem 5] requires the outer function to be convex, univariate, and deterministic, while the inner function can be random but smooth and Vapnik–Chervonenkis-major (VC-major) [22, Section 2.6.4, Exercise 2.6.13]. It remains open whether the VC-major assumption can be dispensed with. Resolving these questions would deepen our understanding of sample-average approximation (SAA) from stochastic programming and empirical risk minimization in machine learning.

In this paper, we give negative answers to all these questions by providing explicit counterexamples. Our negative results hold under natural assumptions, requiring only simple nonsmooth components. They highlight the obstacles that nonsmoothness poses to uniform results. We report the main results in Sections 2.1 and 2.2, followed by a discussion in Section 2.3. All deferred proofs are collected in Section 3.

Notation. We use mostly standard notation. The Borel σ -algebra on Ξ is denoted by $\mathcal{B}(\Xi)$. A random function $f : \Xi \times \mathbb{R}^d \rightarrow \mathbb{R}$ is *Carathéodory* if $f(\xi, \cdot)$ is continuous for each ξ and $f(\cdot, x)$ is (Borel) measurable for each x ; whenever defined, its *gradient* and *Clarke subdifferential* at y are denoted by $\nabla_x f(\xi, y) = \nabla(f(\xi, \cdot))(y)$ and $\partial_x f(\xi, y) = \partial(f(\xi, \cdot))(y)$, respectively. Let $\|\cdot\|$ be the *Euclidean norm* and let $\mathbb{B}(x, r) = \{y \mid \|y - x\| \leq r\}$. For $\nu \in \mathbb{N}$, we write $[\nu]$ for $\{1, \dots, \nu\}$. A function $g : \mathbb{R}^d \rightarrow \mathbb{R}$ is L -Lipschitz on X (or simply L -Lipschitz when $X = \mathbb{R}^d$) if $|g(x) - g(y)| \leq L\|x - y\|$ for all $x, y \in X$; it is L -smooth on X if it is continuously differentiable (C^1) with L -Lipschitz ∇g on X . A set-valued $S : \Xi \rightrightarrows \mathbb{R}^d$ is measurable if $\{\xi \in \Xi \mid S(\xi) \cap C \neq \emptyset\} \in \mathcal{B}(\Xi)$ for every open $C \subset \mathbb{R}^d$; its expectation $\mathbb{E}[S(\boldsymbol{\xi})] \subset \mathbb{R}^d$ is defined in the Aumann sense; see [20, Definition 7.39]. For nonempty sets $C, D \subset \mathbb{R}^d$ and a point $x \in \mathbb{R}^d$, the distance between x and C is $\text{dist}(x, C) = \inf_{z \in C} \|x - z\|$; the excess of C relative to D is $\text{exs}(C; D) = \sup_{z \in C} \text{dist}(z, D)$; and the Hausdorff distance between C and D is $d(C, D) = \max\{\text{exs}(C; D), \text{exs}(D; C)\}$. The $\{0, 1\}$ -indicator function of a set C is defined by $\mathbf{1}_C(x) = 1$ if $x \in C$ and $\mathbf{1}_C(x) = 0$ otherwise. For $\varepsilon \geq 0$, the ε -subdifferential of a convex function $h : \mathbb{R}^d \rightarrow \mathbb{R}$ is defined as $\partial^\varepsilon h(x) = \{s \in \mathbb{R}^d \mid h(y) - h(x) \geq s^\top(y - x) - \varepsilon, \forall y\}$, which coincides with the usual convex subdifferential when $\varepsilon = 0$.

2 Main Results

2.1 Random Lipschitz Functions

We begin with a general negative result for univariate random Lipschitz functions.

Theorem 1 (random Lipschitz). *Let $X = [0, 1] \subset \mathbb{R}$ and ξ^1, ξ^2, \dots be iid random variables on the complete probability space $(\Omega, \mathcal{F}, \mathbb{P})$, each uniformly distributed on $\Xi = [0, 1]$. There exist a Carathéodory $f : \Xi \times \mathbb{R} \rightarrow \mathbb{R}$ and $\delta^\nu \in (0, 1] \downarrow 0$ such that the following hold.*

- (a) *For any $\xi \in \Xi$, the function $f(\xi, \cdot)$ is 1-Lipschitz.*
- (b) *For any $x \in \mathbb{R}$, $|\mathbb{E}[f(\xi, x)]| < \infty$.*
- (c) *For any $\nu \in \mathbb{N}$, the set $D^\nu = \{x \in X \mid f(\xi, \cdot) \text{ is } C^1 \text{ on } \mathbb{B}(x, \delta^\nu), \forall \xi \in \Xi\}$ is nonempty.*
- (d) *\mathbb{P} -almost surely, one has*

$$\liminf_{\nu \rightarrow \infty} \sup_{x \in D^\nu} \inf_{y, \hat{y} \in \mathbb{B}(x, \delta^\nu)} \left| \mathbb{E}[\nabla_x f(\xi, y)] - \frac{1}{\nu} \sum_{i=1}^{\nu} \nabla_x f(\xi^i, \hat{y}) \right| \geq \frac{1}{2}.$$

The definition of the set D^ν requires that, for every $x \in D^\nu$, the function $f(\xi, \cdot)$ is C^1 on $\mathbb{B}(x, \delta^\nu)$ for all $\xi \in \Xi$, so that $\nabla_x f(\xi, \cdot)$ is well defined (hence measurable) and continuous on $\mathbb{B}(x, \delta^\nu)$. This is a stringent requirement, and the fact that we can show $D^\nu \neq \emptyset$ in (c) and restrict to D^ν in (d) only strengthens the theorem. The following corollary is immediate from Theorem 1.

Corollary 1 (subdifferential). *Under the assumptions of Theorem 1, and for $X, \Xi, \{\xi^\nu\}_\nu, f$, and $\{\delta^\nu\}_\nu$ constructed there, the following hold.*

- (a) *For any $\xi \in \Xi$, $\partial_x f(\xi, \cdot)$ is outer semicontinuous with $\text{dist}(0, \partial_x f(\xi, x)) \leq 1$ for all $x \in \mathbb{R}$.*
- (b) *For any $x \in \mathbb{R}$, $\partial_x f(\cdot, x)$ is measurable.*
- (c) *For any sequence $\{r^\nu \in [0, \delta^\nu]\}_\nu$, \mathbb{P} -almost surely, one has*

$$\liminf_{\nu \rightarrow \infty} \sup_{x \in X} \text{exs} \left(\frac{1}{\nu} \sum_{i=1}^{\nu} \partial_x f(\xi^i, x); \bigcup_{y \in \mathbb{B}(x, r^\nu)} \mathbb{E}[\partial_x f(\xi, y)] \right) \geq \frac{1}{2}, \quad (4)$$

$$\liminf_{\nu \rightarrow \infty} \sup_{x \in X} \text{exs} \left(\mathbb{E}[\partial_x f(\xi, x)]; \bigcup_{y \in \mathbb{B}(x, r^\nu)} \frac{1}{\nu} \sum_{i=1}^{\nu} \partial_x f(\xi^i, y) \right) \geq \frac{1}{2}. \quad (5)$$

Proof. The outer semicontinuity, boundedness, and measurability are from standard arguments in the beginning of [19, Section 3]. For measurability, see also [11, Lemma 4]. Part (c) follows from the facts that $\partial_x f(\xi, \cdot) = \{\nabla_x f(\xi, \cdot)\}$ on $\mathbb{B}(x, \delta^\nu)$ for any $x \in D^\nu \subset X$ and (Ω, \mathcal{F}) is \mathbb{P} -complete. \square

In Corollary 1, we show that neither (1) nor (2) can hold when $r = 0$, so (3) is in general impossible for random Lipschitz functions. Consequently, this gives negative answers to the questions posed in [19, Remark 2]. Moreover, it shows that the lower bounds in (4) and (5) remain valid even for $r^\nu \downarrow 0$ shrinking sufficiently fast, so that a trivial “smoothing” technique cannot circumvent our negative result.

It is worth noting that the restriction to the subset $D^\nu \subset X$ in Theorem 1 yields applications far beyond those for the Clarke subdifferential in Corollary 1. For any $\xi \in \Xi$, the function $f(\xi, \cdot)$ is actually smooth on D^ν , and most standard notions of subdifferential coincide there; see, e.g., [14, Chapter 8]. Even for the Clarke subdifferential, we may interchange the sum/expectation and the subdifferential operator, replacing $\frac{1}{\nu} \sum_{i=1}^\nu \partial_x f(\xi^i, x)$ by $\partial(\frac{1}{\nu} \sum_{i=1}^\nu f(\xi^i, \cdot))(x)$ in (4) and $\mathbb{E}[\partial_x f(\xi, x)]$ by $\partial \mathbb{E}[f(\xi, \cdot)](x)$ in (5), while the lower bounds there still hold. Consequently, the failure of uniform laws is not an artifact of the particular subdifferential employed, but reflects a more intrinsic obstruction that extends to a broad class of local approximation concepts.

2.2 Random Convex Functions

Section 2.1 shows that the uniform laws in (1) and (2) with $r = 0$ cannot hold for random Lipschitz functions. One might speculate that Lipschitz functions are somewhat “wild” and the situation would be better for other nonsmooth functions. A natural class to examine next is that of random convex functions, since convexity is widely regarded as favorable. The following uniform law for ε -subdifferentials is classical in the literature.

Theorem 2 (cf. [18, Proposition 3.4] and [20, Theorem 7.56]). *Let nonempty $X \subset \mathbb{R}^d$ be compact, $\Xi \subset \mathbb{R}^m$, and ξ^1, ξ^2, \dots be Ξ -valued iid random variables on the complete probability space $(\Omega, \mathcal{F}, \mathbb{P})$. Let $f : \Xi \times \mathbb{R}^d \rightarrow \mathbb{R}$ be a Carathéodory function such that $f(\xi, \cdot)$ is convex for any $\xi \in \Xi$. Moreover, assume that $|\mathbb{E}[f(\xi, x)]| < \infty$ for any $x \in \mathbb{R}^d$. Then, for any fixed $\varepsilon > 0$, \mathbb{P} -almost surely, one has*

$$\lim_{\nu \rightarrow \infty} \sup_{x \in X} d\left(\partial^\varepsilon \mathbb{E}[f(\xi, \cdot)](x), \partial^\varepsilon \left(\frac{1}{\nu} \sum_{i=1}^\nu f(\xi^i, \cdot)\right)(x)\right) = 0. \quad (6)$$

Theorem 2 differs in several ways from the enlarged uniform laws in (1) and (2). First, there is no Lipschitz assumption in Theorem 2. Second, for a convex h , unlike the r -enlargement $\cup_{y \in \mathbb{B}(x, r)} \partial h(y)$, the ε -subdifferential $\partial^\varepsilon h(x)$ is intrinsically a global notion by definition. Third, the mapping $x \mapsto \partial^\varepsilon h(x)$ is continuous with respect to the Hausdorff metric, whereas the mapping $x \mapsto \cup_{y \in \mathbb{B}(x, r)} \partial h(\xi, y)$ is not. Moreover, as emphasized in [20, p. 382], for random convex function $f : \Xi \times \mathbb{R}^d \rightarrow \mathbb{R}$ and $\xi^1, \dots, \xi^\nu \in \Xi$, the subdifferential sum rule

$$\partial^\varepsilon \left(\frac{1}{\nu} \sum_{i=1}^\nu f(\xi^i, \cdot)\right)(x) = \frac{1}{\nu} \sum_{i=1}^\nu \partial_x^\varepsilon f(\xi^i, x)$$

holds for $\varepsilon = 0$, but fails for $\varepsilon > 0$ and $\nu > 1$.

Therefore, for uniform laws of the form (1) and/or (2) (and, when both are considered, equivalently (6)) for random convex functions, it is natural to ask whether one can take $r = 0$ (or $\varepsilon = 0$ in (6)). We next see that the answer to this question crucially depends on the dimension d of x .

Univariate Case ($d = 1$)

Since the negative result in Theorem 1 for random Lipschitz functions is based on a univariate construction, for random convex functions we likewise first consider the univariate case.

Proposition 1 (univariate random convex). *For a measurable space (Ξ, \mathcal{A}) , let ξ^1, ξ^2, \dots be Ξ -valued iid random variables on the complete probability space $(\Omega, \mathcal{F}, \mathbb{P})$. Let $f : \Xi \times \mathbb{R} \rightarrow \mathbb{R}$ be a Carathéodory function such that $f(\xi, \cdot)$ is convex and $L(\xi)$ -Lipschitz for any $\xi \in \Xi$ with measurable L and $\mathbb{E}[L(\xi)] < \infty$. Assume that $|\mathbb{E}[f(\xi, x)]| < \infty$ for any $x \in \mathbb{R}$. Then, \mathbb{P} -almost surely, one has*

$$\lim_{\nu \rightarrow \infty} \sup_{x \in \mathbb{R}} d\left(\partial \mathbb{E}[f(\xi, \cdot)](x), \frac{1}{\nu} \sum_{i=1}^\nu \partial_x f(\xi^i, x)\right) = 0.$$

In sharp contrast to the univariate random Lipschitz case, convexity allows us to take $r = 0$ in both (1) and (2), even without any compactness assumption. The uniform law in Proposition 1 holds on the entire real line. In spirit, it is similar to the classical Glivenko–Cantelli theorem and, even when restricted to a compact X , already covers cases that were previously unknown. For illustration, let $X = [0, 3]$, and let $\boldsymbol{\xi}^\nu = (\xi_1^\nu, \xi_2^\nu)$ be iid random variables with ξ_1^ν uniformly distributed on $[0, 1]$ and $\xi_2^\nu = 2$ almost surely. Define $f(\xi, x) = \max\{x - \xi_1, 0\} + \max\{\xi_2 - x, 0\}$. It is easy to see that $x \mapsto \mathbb{E}[f(\boldsymbol{\xi}, x)]$ is not differentiable (it has a kink at $x = 2$), and a computation similar to [12, Example 3.4] shows that the range of $\partial_x f$ is not essentially separable. Moreover, one can show that f cannot be written in the univariate convex-composite form considered in [15, Section 4]. Hence, a subdifferential uniform law for f cannot be deduced from the results in [17], [21], or [15]. However, Proposition 1 applies.

For the univariate case ($d = 1$), we thus have obtained reasonably satisfactory positive results. What about higher dimensions?

General Case ($d \geq 2$)

Perhaps surprisingly, the situation changes drastically as soon as $d = 2$.

Theorem 3 (random convex). *Let $X = \{0\} \times [0, 1] \subset \mathbb{R}^2$, and $\boldsymbol{\xi}^1, \boldsymbol{\xi}^2, \dots$ be iid random variables on the complete probability space $(\Omega, \mathcal{F}, \mathbb{P})$, each uniformly distributed on $\Xi = [0, 1]$. There exist a Carathéodory $g : \Xi \times \mathbb{R}^2 \rightarrow \mathbb{R}$ and $\delta^\nu \in (0, 1] \downarrow 0$ such that for*

$$(\xi, x) \mapsto f(\xi, x) = \max\{g(\xi, x), 0\} + 35\|x\|^2,$$

the following hold.

- (a) *For any $\xi \in \Xi$, $g(\xi, \cdot)$ is 70-Lipschitz, 70-smooth, and $f(\xi, \cdot)$ is convex, 140-Lipschitz on X .*
- (b) *For any $x \in \mathbb{R}^2$, $|\mathbb{E}[f(\boldsymbol{\xi}, x)]| < \infty$.*
- (c) *For any $\nu \in \mathbb{N}$, the set $D^\nu = \{x \in X \mid f(\xi, \cdot) \text{ is } C^1 \text{ on } \mathbb{B}(x, \delta^\nu), \forall \xi \in \Xi\}$ is nonempty.*
- (d) *\mathbb{P} -almost surely, one has*

$$\liminf_{\nu \rightarrow \infty} \sup_{x \in D^\nu} \inf_{y, \hat{y} \in \mathbb{B}(x, \delta^\nu)} \left\| \mathbb{E}[\nabla_x f(\boldsymbol{\xi}, y)] - \frac{1}{\nu} \sum_{i=1}^{\nu} \nabla_x f(\boldsymbol{\xi}^i, \hat{y}) \right\| \geq \frac{1}{2}.$$

In Theorem 3, we show that even in dimension $d = 2$ there exists a random convex function $f : \Xi \times \mathbb{R}^2 \rightarrow \mathbb{R}$ for which a negative result analogous to Theorem 1 holds. Therefore, in the same vein as Corollary 1, one cannot take $r = 0$ in (1) and (2); equivalently, a uniform law with $\varepsilon = 0$ in (6) cannot hold. The same result holds for any $d \geq 2$ simply by padding with zero coordinates.

However, Theorem 3 says more than this, since the random convex f there is highly structured. Let us make the following observations:

- In contrast to the construction used to prove Theorem 1 (see Section 3.2), which involves countably infinitely many pieces¹ (see Figure 1), the function f in Theorem 3 is a random

¹Here, by the number of “pieces,” we mean the smallest number of connected regions partitioning the domain such that the function is smooth on the interior of each region.

convex function with only two pieces (see Figure 2). Indeed, the only source of nonsmoothness in the construction is the function $t \mapsto \max\{t, 0\}$. Moreover, by [20, Theorem 7.52], we have

$$\lim_{\nu \rightarrow \infty} \sup_{x \in X} \left\| \mathbb{E}[\nabla_x g(\xi, x)] - \frac{1}{\nu} \sum_{i=1}^{\nu} \nabla_x g(\xi^i, x) \right\| = 0$$

for the function g emerging from Theorem 3, \mathbb{P} -almost surely, with X being any nonempty compact set. Thus, each smooth component appearing in the definition of f satisfies a uniform law on its own. However, when they are combined through this (arguably) simplest nonsmooth operation $t \mapsto \max\{t, 0\}$, the uniform law fails. This phenomenon is driven by the oscillating boundary between the connected sets $\{x \mid g(\xi, x) > 0\}$ and $\{x \mid g(\xi, x) < 0\}$, whose intersection with $X = \{0\} \times [0, 1]$ gives rise to countably infinitely many kinks in the univariate function $t \mapsto \max\{g(\xi, (0, t)), 0\}$.

- The negative results in Theorem 3 hold for any function class that contains f . An interesting example is a subclass of the random convex-composite functions studied in [15, Section 4]. For $\kappa < \infty$, consider a convex, *piecewise affine*, κ -Lipschitz $h : \mathbb{R} \rightarrow \mathbb{R}$ and a Carathéodory $c : \Xi \times \mathbb{R}^d \rightarrow \mathbb{R}$ such that $c(\xi, \cdot)$ is κ -Lipschitz and κ -smooth for each ξ . Since κ is finite and h is piecewise affine, the constants in [15, Assumptions C.1–C.3] are *distribution-free*. By [15, Theorem 5], for large ν , with (outer) probability at least $1 - \delta$, one has

$$\sup_{x \in B} d \left(\partial \mathbb{E}[(h \circ c)(\xi, \cdot)](x), \frac{1}{\nu} \sum_{i=1}^{\nu} \partial_x (h \circ c)(\xi^i, x) \right) \leq O \left(\sqrt{\frac{d + \text{VC}(\mathcal{C}) \log \nu + \log(1/\delta)}{\nu}} \right),$$

where B is an open Euclidean ball, $\mathcal{C} = \{\{\xi \mid c(\xi, x) \geq t\} \mid x \in \mathbb{R}^d, t \in \mathbb{R}\}$, and $\text{VC}(\mathcal{C})$ denotes the VC dimension of \mathcal{C} ; see [22, Section 2.6.1]. This upper bound is also *distribution-free*, but is non-trivial only if $\text{VC}(\mathcal{C}) < \infty$; in other words, the function class $\{c(\cdot, x) \mid x \in \mathbb{R}^d\}$ is VC-major; see [22, Section 2.6.4, Exercise 2.6.13]. The VC-type assumptions are usually employed to obtain *distribution-free* results [22, Section 2.8.1] and may not be necessary for universal Glivenko–Cantelli-type results, let alone $\mathbb{P} \circ (\xi^1)^{-1}$ -Glivenko–Cantelli-type results such as [20, Theorem 7.52]. Hence, it is natural to ask whether the VC-major assumption here can be dispensed with when uniformity over the underlying distribution and an explicit convergence rate are not required. Consider the convex piecewise affine $h = \max\{\cdot, 0\}$ and let $c = g$, where g is as in Theorem 3. We can write the function f in Theorem 3 as

$$f(\xi, x) = h(c(\xi, x)) + 35\|x\|^2.$$

Ignoring the deterministic term $x \mapsto 35\|x\|^2$, this is exactly a random convex-composite function of the above form with $\kappa = 70$ but $\text{VC}(\mathcal{C}) = \infty$.² When X in Theorem 3 is a subset of B , we see that, even without requiring uniformity over distributions, the VC-major assumption cannot, in general, be dropped without causing the uniform law to fail.

- Another example is median regression. Let $\xi^\nu = (\xi_1^\nu, \xi_2^\nu)$ be iid $\Xi \times [-1, 1]$ -valued random variables with $\Xi \subset \mathbb{R}^m$, and let $\psi : \Xi \times \mathbb{R}^d \rightarrow \mathbb{R}$ be a Carathéodory function such that $\psi(\xi_1, \cdot)$

²To see it, let $\text{bit}_k : [0, 1] \rightarrow \{0, 1\}$ and $g : \Xi \times \mathbb{R}^2 \rightarrow \mathbb{R}$ be defined in Sections 3.1 and 3.4. When $x = (0, 1/k) \in X$ and $t = 0$, one has $\{\xi \mid g(\xi, x) \geq t\} = \{\xi \mid \text{bit}_k(\xi) = 1\}$. To show $\text{VC}(\mathcal{C}) = \infty$, it suffices to show that, for any $n \in \mathbb{N}$, there exist $\xi^1, \dots, \xi^n \in \Xi$ such that, for any $b \in \{0, 1\}^n$, one has $\text{bit}_k(\xi^i) = b_i$ for each $i \in [n]$ and some $k \in \mathbb{N}$. Let $\{b^k \mid k \in [2^n]\} = \{0, 1\}^n$ cycle through all 2^n patterns. A valid construction is $\xi^i = \sum_{k=1}^{2^n} 2^{-k} b_i^k$ for each $i \in [n]$.

is Lipschitz and smooth for each $\xi_1 \in \Xi$. Median regression corresponds to the random function $f(\xi, x) = |\xi_2 - \psi(\xi_1, x)|$. It is easy to see that the construction in Theorem 3 can be rewritten in this form, since $f(\xi, x) = \max\{2\xi_2 - 2\psi(\xi_1, x), 0\} - \xi_2 + \psi(\xi_1, x)$ with the last term being handled by [20, Theorem 7.52]. Therefore, in general, uniform laws fail for the subdifferential of f when $d > 1$.

It is straightforward to verify that the conclusion of Corollary 1 still holds under the assumptions of Theorem 3, and we omit the details for brevity. We end this subsection with a negative result for ε -subdifferentials, complementing the positive result in Theorem 2 and addressing the case of rapidly shrinking $\varepsilon^\nu \downarrow 0$.

Corollary 2 (ε -subdifferential). *Under the assumptions of Theorem 3, let X , Ξ , $\{\xi^\nu\}_\nu$, f , and $\{\delta^\nu\}_\nu$ be constructed there. For any sequence $\{\varepsilon^\nu \in [0, (\delta^\nu)^2]\}_\nu$, \mathbb{P} -almost surely, one has*

$$\liminf_{\nu \rightarrow \infty} \sup_{x \in X} \text{exs} \left(\partial^{\varepsilon^\nu} \left(\frac{1}{\nu} \sum_{i=1}^{\nu} f(\xi^i, \cdot) \right)(x); \partial^{\varepsilon^\nu} \mathbb{E}[f(\xi, \cdot)](x) \right) \geq \frac{1}{2},$$

$$\liminf_{\nu \rightarrow \infty} \sup_{x \in X} \left(\partial^{\varepsilon^\nu} \mathbb{E}[f(\xi, \cdot)](x); \partial^{\varepsilon^\nu} \left(\frac{1}{\nu} \sum_{i=1}^{\nu} f(\xi^i, \cdot) \right)(x) \right) \geq \frac{1}{2}.$$

Proof. For any $x^\nu \in D^\nu$ in Theorem 3, let $g_1^\nu \in \partial^{\varepsilon^\nu} \mathbb{E}[f(\xi, \cdot)](x^\nu)$ and $g_2^\nu \in \partial^{\varepsilon^\nu} \left(\frac{1}{\nu} \sum_{i=1}^{\nu} f(\xi^i, \cdot) \right)(x^\nu)$. By Brøndsted–Rockafellar [2, p. 608] and $\varepsilon^\nu \in [0, (\delta^\nu)^2]$, we get $\|g_1^\nu - \mathbb{E}[\nabla_x f(\xi, y_1^\nu)]\| \leq \delta^\nu$ and $\|g_2^\nu - \frac{1}{\nu} \sum_{i=1}^{\nu} \nabla_x f(\xi^i, y_2^\nu)\| \leq \delta^\nu$, where $y_1^\nu, y_2^\nu \in \mathbb{B}(x^\nu, \delta^\nu)$. Then, $\|g_1^\nu - g_2^\nu\| \geq \|\mathbb{E}[\nabla_x f(\xi, y_1^\nu)] - \frac{1}{\nu} \sum_{i=1}^{\nu} \nabla_x f(\xi^i, y_2^\nu)\| - 2\delta^\nu$. The proof completes by invoking Theorem 3 and $\delta^\nu \downarrow 0$. \square

2.3 Discussion

The previous subsections show that, for certain seemingly benign nonsmooth functions, obtaining a uniform LLN for subdifferentials is impossible. These negative results rule out a wide class of functions, but also shed light on where positive results might still exist. All our constructions involve functions that oscillate near some point; thus, at least in order to circumvent our counterexamples, some notion of tameness could be helpful [7].

Nevertheless, by considering weaker notions of convergence for subdifferentials, positive results can be obtained for a broad range of functions. One such notion is graphical convergence; see [12] for general results on graphical LLNs, [4] for convergence rates for weakly convex functions, and also [16] for results on monotone operators. Indeed, under suitable boundedness and joint measurability assumptions, it is shown in [12, Proposition 4.3] that the uniform law for *enlarged* subdifferentials in [19, Theorem 5] is equivalent to a graphical LLN on a compact set. Thus, our negative results can be interpreted as a *separation* between uniform and graphical LLNs for set-valued mappings.

We end with a concrete implication of this separation. For a random convex f , let x_0 be any ε -stationary point of $x \mapsto \frac{1}{\nu} \sum_{i=1}^{\nu} f(\xi^i, x)$; i.e., $\text{dist}(0, \partial(\frac{1}{\nu} \sum_{i=1}^{\nu} f(\xi^i, \cdot))(x_0)) \leq \varepsilon$. If the uniform LLN in (3) holds and ν is sufficiently large, then we can conclude that x_0 is a 2ε -stationary point of $x \mapsto \mathbb{E}[f(\xi, x)]$. In contrast, when only a graphical LLN holds, we can guarantee only the existence of some point x'_0 near x_0 such that x'_0 is 2ε -stationary for $x \mapsto \mathbb{E}[f(\xi, x)]$, and locating such an x'_0 may be difficult even though its existence is assured.

3 Proofs

We now present the proofs of our main results, beginning with a random bits gadget.

3.1 A Probabilistic Gadget

For any $\xi \in [0, 1]$ and $k \in \mathbb{N}$, let

$$\text{bit}_k(\xi) = \lfloor 2^k \xi \rfloor - 2 \lfloor 2^{k-1} \xi \rfloor \in \{0, 1\},$$

which can be understood as the k th binary digit of $\xi - \lfloor \xi \rfloor$ using the terminating-zeros convention for dyadic rationals. For $\xi \in [0, 1)$, we have the binary expansion

$$\xi = \sum_{k=1}^{\infty} 2^{-k} \text{bit}_k(\xi).$$

Let ξ^1, ξ^2, \dots be iid random variables on the complete probability space $(\Omega, \mathcal{F}, \mathbb{P})$, each with uniform distribution on $\Xi = [0, 1]$. It is classical that the random variables $\{\text{bit}_k(\xi^i)\}_{i,k \in \mathbb{N}}$ are iid Bernoulli random variables with parameter $\frac{1}{2}$; see, e.g., [23, Section 4.6] and [8, Lemma 2.20].

Lemma 1. *For any $k \in \mathbb{N}$, the function $\xi \mapsto \text{bit}_k(\xi)$ is $\mathcal{B}(\Xi)$ -measurable.*

Proof. Simply observe that $t \mapsto \lfloor t \rfloor$ is $\mathcal{B}(\Xi)$ -measurable. \square

Lemma 2. *There exists an \mathcal{F} -measurable $E \subset \Omega$ with $\mathbb{P}(E) = 1$ such that for any $\omega \in E$ and any $\nu \geq \bar{\nu}(\omega)$, one has*

$$\text{bit}_{\mathbf{k}^\nu(\omega)}(\xi^1(\omega)) = \dots = \text{bit}_{\mathbf{k}^\nu(\omega)}(\xi^\nu(\omega)) = 1$$

for some $\bar{\nu}(\omega) \in \mathbb{N}$ and $\mathbf{k}^\nu(\omega) \leq \lceil 2^{\nu+1} \log(\nu+1) \rceil$.

Proof. Let $K^\nu = \lceil 2^{\nu+1} \log(\nu+1) \rceil$. For any $\nu, k \in \mathbb{N}$, define sets

$$E_k^\nu = \bigcap_{i \in [\nu]} \{\text{bit}_k(\xi^i) = 1\}, \quad E = \bigcup_{\bar{\nu} \geq 1} \bigcap_{\nu \geq \bar{\nu}} \bigcup_{k \in [K^\nu]} E_k^\nu,$$

which are \mathcal{F} -measurable from Lemma 1. By independence of $\{\text{bit}_k(\xi^i)\}_{i \in [\nu], k \in \mathbb{N}}$, we have

$$\begin{aligned} \mathbb{P}\left(\bigcap_{k \in [K^\nu]} (E_k^\nu)^c\right) &= \prod_{k \in [K^\nu]} \mathbb{P}((E_k^\nu)^c) = \prod_{k \in [K^\nu]} \left(1 - \prod_{i=1}^\nu \mathbb{P}(\{\text{bit}_k(\xi^i) = 1\})\right) \\ &= (1 - 2^{-\nu})^{K^\nu} \leq e^{-2^{-\nu} K^\nu} \leq \frac{1}{(\nu+1)^2}. \end{aligned}$$

Then, $\sum_{\nu=1}^{\infty} \mathbb{P}(\bigcap_{k \in [K^\nu]} (E_k^\nu)^c) < \infty$. By Borel–Cantelli, we have

$$\mathbb{P}\left(\bigcap_{\bar{\nu} \geq 1} \bigcup_{\nu \geq \bar{\nu}} \bigcap_{k \in [K^\nu]} (E_k^\nu)^c\right) = \mathbb{P}(E^c) = 0.$$

For any $\nu \in \mathbb{N}$, define $\mathbf{k}^\nu(\omega) = 1$ if $\omega \notin \bigcup_{k \in [K^\nu]} E_k^\nu$ and $\bar{\nu}(\omega) = 1$ if $\omega \notin E$. Otherwise, set

$$\mathbf{k}^\nu(\omega) = \min\{k \in [K^\nu] \mid \omega \in E_k^\nu\}, \quad \bar{\nu}(\omega) = \min\{\bar{\nu} \geq 1 \mid \omega \in \bigcap_{\nu \geq \bar{\nu}} \bigcup_{k \in [K^\nu]} E_k^\nu\},$$

which are \mathcal{F} -measurable and complete the proof by construction. \square

3.2 Proof of Theorem 1

Intuition. Our proof is inspired by a construction that refutes one-sided uniform consistency in the setting of epigraphical LLNs; see [1, Example 4.1]. The key idea is to accumulate countably infinitely many affine pieces near the origin, with the slope of each piece determined by a random bit, $\text{bit}_k(\xi^i) \in \{0, 1\}$. Then, with the help of Lemma 2, we can almost surely find a random point in D^ν at which the functions $\{f(\xi^i, \cdot) \mid i \in [\nu]\}$ are all C^1 with gradient consistently equal to 1, thereby exhibiting a constant gap between the empirical and the average cases. Our proof is divided into five steps.

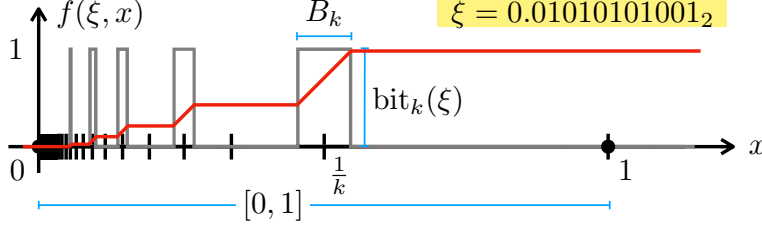


Figure 1: The function $f(\xi, \cdot)$ in Section 3.2 when $\xi = 0.01010101001_2$.

Step 1. Define open sets $B_k = \{x \in \mathbb{R} \mid |x - \frac{1}{k}| < r_k\}$ with $r_k = \frac{1}{4k^2}$ for $k \in \mathbb{N}$. Let

$$g(\xi, x) = \sum_{k=1}^{\infty} \mathbf{1}_{B_k}(x) \text{bit}_k(\xi).$$

Since $\frac{1}{2}(\frac{1}{k} - \frac{1}{k+1}) \geq \frac{1}{4k^2} = r_k$, the sets $\{B_k\}_k$ are disjoint, so that $g(\xi, x) \in \{0, 1\}$. Let

$$f(\xi, x) = \int_0^{\max\{x, 0\}} g(\xi, t) dt,$$

which is well-defined as explained below. Define $\delta^\nu = (8 \lceil 2^{\nu+1} \log(\nu+1) \rceil^2)^{-1}$.

Step 2. The $\{0, 1\}$ -valued function $g(\xi, \cdot)$ is a sum of measurable indicator functions, then also measurable. Since $g(\xi, \cdot)$ is bounded and has only countably many points of discontinuity, the Lebesgue criterion implies that $g(\xi, \cdot)$ is Riemann integrable on compact sets, and hence f is well defined. For any $x, y \in \mathbb{R}$, one has $|f(\xi, x) - f(\xi, y)| \leq |\int_{\min\{x, y\}}^{\max\{x, y\}} g(\xi, t) dt| \leq |x - y|$, hence $f(\xi, \cdot)$ is 1-Lipschitz continuous. From Lemma 1, the function $\xi \mapsto \text{bit}_k(\xi)$ is measurable. By Tonelli, we have $f(\xi, x) = \sum_{k=1}^{\infty} \text{bit}_k(\xi) \int_0^{\max\{x, 0\}} \mathbf{1}_{B_k} dt < \infty$, hence $f(\cdot, x)$ is also measurable. Therefore, f is Carathéodory. Moreover, for any $x \in \mathbb{R}$, $|\mathbb{E}[f(\xi, x)]| \leq \int_0^{\max\{x, 0\}} 1 dt < \infty$.

Step 3. Let $p_k = \frac{1}{k}$ and $\Delta_k = \frac{r_k}{2} = \frac{1}{8k^2}$. For any $k \in \mathbb{N}$, $\xi \in \Xi$, and $y \in \mathbb{B}(p_k, \Delta_k) \subset B_k$, one has

$$g(\xi, y) = \text{bit}_k(\xi), \quad f(\xi, y) = \text{bit}_k(\xi)(y - p_k + \Delta_k) + \int_0^{p_k - \Delta_k} g(\xi, t) dt.$$

Hence, the function, $f(\xi, \cdot)$ is smooth on $\cup_k \mathbb{B}(p_k, \Delta_k)$. Moreover, for any $k \in \mathbb{N}$ and $y \in \mathbb{B}(p_k, \Delta_k)$, we can write $\nabla_x f(\xi, y) = g(\xi, y) = \text{bit}_k(\xi)$.

Step 4. Let \mathcal{F} -measurable $E \subset \Omega$ with $\mathbb{P}(E) = 1$, $\bar{\nu}$, and $\{\mathbf{k}^\nu\}_\nu$ be defined in Lemma 2. For any $\omega \in E$ and $\nu \geq \bar{\nu}(\omega)$, one has $\text{bit}_{\mathbf{k}^\nu(\omega)}(\xi^i(\omega)) = 1$ for all $i \in [\nu]$ and some $\mathbf{k}^\nu(\omega) \leq \lceil 2^{\nu+1} \log(\nu+1) \rceil$. Let $\mathbf{p}^\nu(\omega) = p_{\mathbf{k}^\nu(\omega)} \in \{p_k\}_k$. For any $y \in \mathbb{B}(\mathbf{p}^\nu(\omega), \Delta_{\mathbf{k}^\nu(\omega)})$, we have

$$\frac{1}{\nu} \sum_{i=1}^{\nu} \nabla_x f(\xi^i(\omega), y) = \frac{1}{\nu} \sum_{i=1}^{\nu} \text{bit}_{\mathbf{k}^\nu(\omega)}(\xi^i(\omega)) = 1.$$

Meanwhile, for any $k \in \mathbb{N}$ and $y \in \mathbb{B}(p_k, \Delta_k)$, deterministically, we have

$$\mathbb{E}[\nabla_x f(\xi, y)] = \mathbb{E}_\xi[\text{bit}_k(\xi)] = \frac{1}{2}.$$

Step 5. Note that

$$0 \leq \delta^\nu = (8 \lceil 2^{\nu+1} \log(\nu+1) \rceil^2)^{-1} = \Delta_{\lceil 2^{\nu+1} \log(\nu+1) \rceil} \leq \Delta_{\mathbf{k}^\nu(\omega)}.$$

Therefore, for any $\omega \in E$, we conclude that $\mathbf{p}^\nu(\omega) \in \{p_k \mid 1 \leq k \leq \lceil 2^{\nu+1} \log(\nu+1) \rceil\} \subset D^\nu$ and

$$\begin{aligned} & \liminf_{\nu \rightarrow \infty} \sup_{x \in D^\nu} \inf_{y, \hat{y} \in \mathbb{B}(x, \delta^\nu)} \left| \mathbb{E}[\nabla_x f(\boldsymbol{\xi}, y)] - \frac{1}{\nu} \sum_{i=1}^\nu \nabla_x f(\boldsymbol{\xi}^i(\omega), \hat{y}) \right| \\ & \geq \inf_{\nu \geq \bar{\nu}(\omega)} \inf_{y, \hat{y} \in \mathbb{B}(\mathbf{p}^\nu(\omega), \delta^\nu)} \left| \mathbb{E}[\nabla_x f(\boldsymbol{\xi}, y)] - \frac{1}{\nu} \sum_{i=1}^\nu \nabla_x f(\boldsymbol{\xi}^i(\omega), \hat{y}) \right| = |\tfrac{1}{2} - 1| = \tfrac{1}{2}, \end{aligned}$$

which completes the proof, since (Ω, \mathcal{F}) is \mathbb{P} -complete.

3.3 Proof of Proposition 1

Since $f(\xi, \cdot)$ is real-valued and convex, by [13, Theorem 23.1], the directional derivative of $f(\xi, \cdot)$ exists at any point x and in every direction w , and we denote it by $f'_x(\xi, x; w)$. Using [20, Theorem 7.46], the function $x \mapsto \mathbb{E}[f(\boldsymbol{\xi}, x)]$ is convex and directionally differentiable with \mathcal{A} -measurable $f'_x(\cdot, x; w)$ and $(\mathbb{E}[f(\boldsymbol{\xi}, \cdot)])'(x; w) = \mathbb{E}[f'_x(\boldsymbol{\xi}, x; w)]$. From [6, Theorems V.3.3.8, V.3.3.3], [13, Theorem 23.2], and the non-emptiness of compact $\partial_x f(\xi, x)$, we can write the Hausdorff distance as follows

$$d\left(\partial \mathbb{E}[f(\boldsymbol{\xi}, \cdot)](x), \frac{1}{\nu} \sum_{i=1}^\nu \partial_x f(\boldsymbol{\xi}^i, x)\right) = \sup_{|w|=1} \left| \mathbb{E}[f'_x(\boldsymbol{\xi}, x; w)] - \frac{1}{\nu} \sum_{i=1}^\nu f'_x(\boldsymbol{\xi}^i, x; w) \right|.$$

Let $P = \mathbb{P} \circ (\boldsymbol{\xi}^1)^{-1}$. Hence, it suffices to show that, for any fixed $w \in \{-1, 1\}$, the class of \mathcal{A} -measurable functions $\mathcal{S} = \{f'_x(\cdot, x; w) \mid x \in \mathbb{R}\}$ is P -Glivenko–Cantelli; i.e.,

$$\mathbb{P} \left(\lim_{\nu \rightarrow \infty} \sup_{x \in \mathbb{R}} \left| \mathbb{E}[f'_x(\boldsymbol{\xi}, x; w)] - \frac{1}{\nu} \sum_{i=1}^\nu f'_x(\boldsymbol{\xi}^i, x; w) \right| = 0 \right) = 1, \quad (7)$$

which is well defined as explained below. For simplicity, assume $w = 1$. Using the fact that $f(\xi, \cdot)$ is $L(\xi)$ -Lipschitz, we have $|f'_x(\xi, x; w)| \leq L(\xi)$ for any $x \in \mathbb{R}$. Because $\mathbb{E}[L(\boldsymbol{\xi})] < \infty$, the class \mathcal{S} admits an integrable envelope L . From [20, Theorem 7.42] and [13, Theorem 24.1], the function $t \mapsto f'_x(\xi, t; w)$ is non-decreasing and right-continuous (when $w = -1$, it is non-increasing and left-continuous), thereby \mathcal{S} is pointwise measurable [22, Example 2.3.4] and hence P -measurable [22, Definition 2.3.3] with (7) being well defined. With monotonicity and boundedness of $f'_x(\xi, \cdot; w)$, set

$$l(\xi) = \lim_{k \rightarrow \infty} f'_x(\xi, -k; w), \quad u(\xi) = \lim_{k \rightarrow \infty} f'_x(\xi, k; w),$$

where $l, u : \Xi \rightarrow \mathbb{R}$ are \mathcal{A} -measurable with $\mathbb{E}[u(\boldsymbol{\xi})] - \mathbb{E}[l(\boldsymbol{\xi})] \leq 2\mathbb{E}[L(\boldsymbol{\xi})]$. If $\mathbb{E}[u(\boldsymbol{\xi})] = \mathbb{E}[l(\boldsymbol{\xi})]$, then $f'_x(\xi, \cdot; w)$ is almost surely constant and the claim is trivial. Assume $\mathbb{E}[(u - l)(\boldsymbol{\xi})] > 0$, and fix any $0 < \varepsilon < \frac{2}{3}\mathbb{E}[(u - l)(\boldsymbol{\xi})]$. For each $n \in \mathbb{N}$, define

$$t_n = \min \{t \in \mathbb{R} \mid \mathbb{E}[f'_x(\boldsymbol{\xi}, t; w)] \geq \min\{\mathbb{E}[u(\boldsymbol{\xi})] - \tfrac{\varepsilon}{2}, \mathbb{E}[l(\boldsymbol{\xi})] + n\varepsilon\}\},$$

where the attainment is from the right-continuity of $t \mapsto \mathbb{E}[f'_x(\boldsymbol{\xi}, t; w)]$. Let $N \in \mathbb{N}$ be the first number such that $\mathbb{E}[f'_x(\boldsymbol{\xi}, t_N; w)] \geq \mathbb{E}[u(\boldsymbol{\xi})] - \frac{\varepsilon}{2}$. By construction, we have $N \leq 1 + 2\mathbb{E}[L(\boldsymbol{\xi})]/\varepsilon < \infty$ and $\mathbb{E}[|u(\boldsymbol{\xi}) - f'_x(\boldsymbol{\xi}, t_N; w)|] \leq \varepsilon$. By [13, Theorem 24.1], we get $-f'_x(\xi, t_n; -w) = \lim_{s \uparrow t_n} f'_x(\xi, s; w)$. Using monotone convergence theorem and the definition of t_n , for any $1 \leq n < N$, one has

$$\mathbb{E}[-f'_x(\boldsymbol{\xi}, t_n; -w)] = \lim_{s \uparrow t_n} \mathbb{E}[f'_x(\boldsymbol{\xi}, s; w)] \leq \mathbb{E}[l(\boldsymbol{\xi})] + n\varepsilon,$$

which yields that $\mathbb{E}[| -f'_x(\xi, t_{n+1}; -w) - f'_x(\xi, t_n; w) |] \leq \varepsilon$ and $\mathbb{E}[| -f'_x(\xi, t_1; -w) - l(\xi) |] \leq \varepsilon$. Then, the following sets of functions

$$\begin{aligned} [l, -f'_x(\cdot, t_1; -w)] &= \{\phi \in \mathcal{S} \mid l \leq \phi \leq -f'_x(\cdot, t_1; -w)\}, \\ [f'_x(\cdot, t_n; w), -f'_x(\cdot, t_{n+1}; -w)] &= \{\phi \in \mathcal{S} \mid f'_x(\cdot, t_n; w) \leq \phi \leq -f'_x(\cdot, t_{n+1}; -w)\}, \quad \forall 1 \leq n < N, \\ [f'_x(\cdot, t_N; w), u] &= \{\phi \in \mathcal{S} \mid f'_x(\cdot, t_N; w) \leq \phi \leq u\}, \end{aligned}$$

are ε -brackets of \mathcal{S} in $L_1(P)$ -(semi)norm. Meanwhile, for any $x \in \mathbb{R}$, one has

$$f'_x(\cdot, x; w) \in \bigcup_{n < N} [f'_x(\cdot, t_n; w), -f'_x(\cdot, t_{n+1}; -w)] \cup [l, -f'_x(\cdot, t_1; -w)] \cup [f'_x(\cdot, t_N; w), u].$$

Hence, the bracketing number $\mathcal{N}_{[]}(\varepsilon, \mathcal{S}, L^1(P))$ is finite. Invoking [22, Theorem 2.4.1], we conclude that the class \mathcal{S} is P -Glivenko–Cantelli. Applying the same argument to $w = -1$, with simple adjustments for nonincreasingness and left-continuity, and then taking the maximum over $w \in \{-1, 1\}$ completes the proof.

3.4 Proof of Theorem 3

Intuition. When f is random convex, the “square-wave” function g used in Section 3.2 cannot be extended in a straightforward way, since convex subdifferentials must be monotone. The remedy is to introduce an orthogonal direction and control the activation of the “pieces” by a smoothed “square-wave” function. This smoothed function must have decreasing magnitude near the origin, due to the assumption about Lipschitz gradients, hence cannot create a constant gap on its own. However, we can exploit the tiny oscillations in its function values near the origin, combined with the nonsmooth map $t \mapsto \max\{t, 0\}$ and the orthogonal direction with a constant slope, to amplify the “signal” coming from the random bits. The proof again proceeds in five steps.

Step 1. Let $\rho : \mathbb{R} \rightarrow [0, 1]$ be a twice continuously differentiable function such that

$$\text{supp}(\rho) = [-1, 1], \quad \rho([- \tfrac{1}{2}, \tfrac{1}{2}]) = 1, \quad \rho'([- \tfrac{1}{2}, \tfrac{1}{2}]) = 0, \quad |\rho'(\cdot)| \leq 30, \quad |\rho''(\cdot)| \leq 30,$$

where $\text{supp}(\rho)$ is the *support* of the function ρ . One possible choice is

$$\rho(t) = \begin{cases} 0 & \text{if } |t| \geq 1, \\ 1 - \theta(2|t| - 1) & \text{if } |t| \in (\tfrac{1}{2}, 1), \\ 1 & \text{if } |t| \leq \tfrac{1}{2}, \end{cases}$$

where $\theta(t) = 6t^5 - 15t^4 + 10t^3$. For each $k \in \mathbb{N}$, define $\psi_k : \mathbb{R} \rightarrow [0, 1]$ as

$$\psi_k(t) = \eta_k \rho\left(\frac{t - \frac{1}{k}}{r_k}\right),$$

where $r_k = \frac{1}{8k^2}$, $\eta_k = r_k^2 = \frac{1}{64k^4}$. Correspondingly, one has

$$\text{supp}(\psi_k) = \mathbb{B}(\tfrac{1}{k}, r_k), \quad \psi_k(\mathbb{B}(\tfrac{1}{k}, \tfrac{r_k}{2})) = \eta_k, \quad \psi'_k(\mathbb{B}(\tfrac{1}{k}, \tfrac{r_k}{2})) = 0, \quad |\psi'_k(\cdot)| \leq 30, \quad |\psi''_k(\cdot)| \leq 30. \quad (8)$$

Define a random smooth function $g : \Xi \times \mathbb{R}^2 \rightarrow \mathbb{R}$ as

$$g(\xi, x) = x_1 + \sum_{k=1}^{\infty} \psi_k(x_2)(2 \text{bit}_k(\xi) - 1),$$

which is finite valued as explained below. Let $h = \max\{\cdot, 0\}$ and $f : \Xi \times \mathbb{R}^2 \rightarrow \mathbb{R}$ be

$$f(\xi, x) = h(g(\xi, x)) + 35\|x\|^2 = \max\{g(\xi, x), 0\} + 35\|x\|^2.$$

Define $\delta^\nu = (2240 \lceil 2^{\nu+1} \log(\nu+1) \rceil^4)^{-1}$.

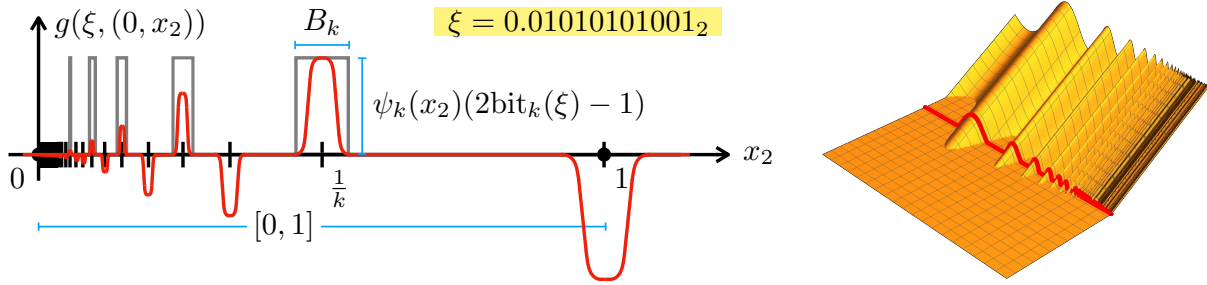


Figure 2: The functions $g(\xi, (0, \cdot))$ and $\max\{g(\xi, \cdot), 0\}$ in Section 3.4 when $\xi = 0.01010101001_2$.

Step 2. Let $C = 70$. Since $\frac{1}{2}(\frac{1}{k} - \frac{1}{k+1}) > \frac{1}{8k^2} = r_k$ for all $k \in \mathbb{N}$, the sets $\{\text{supp}(\psi_k)\}_k$ are disjoint, hence, for any $t \in \mathbb{R}$, at most one element from $\{\psi_k(t)\}_k$ can be positive. For $x = (x_1, x_2), y = (y_1, y_2)$, at most two (say, with indices k_1, k_2) in $\{\psi_k(x_2)\}_k \cup \{\psi_k(y_2)\}_k$ can be positive, which gives

$$\|g(\xi, x) - g(\xi, y)\| \leq |x_1 - y_1| + \sum_{j \in \{k_1, k_2\}} |\psi_j(x_2) - \psi_j(y_2)| \leq 61\|x - y\|.$$

Then, $g(\xi, \cdot)$ is C -Lipschitz continuous for any $\xi \in \Xi$ with

$$\nabla_x g(\xi, x) = \left[\begin{array}{c} 1 \\ \sum_{k=1}^{\infty} \psi'_k(x_2)(2 \text{bit}_k(\xi) - 1) \end{array} \right].$$

Similarly, one has $\|\nabla_x g(\xi, x) - \nabla_x g(\xi, y)\| \leq C\|x - y\|$, so that the function $g(\xi, \cdot)$ is C -smooth. Hence, $f(\xi, \cdot)$ is convex and $2C$ -Lipschitz on X . From Lemma 1, the functions $\xi \mapsto \text{bit}_k(\xi)$, $\xi \mapsto g(\xi, x)$, and then $\xi \mapsto f(\xi, x)$ are measurable. Therefore, f and g are Carathéodory. Moreover, for any $x \in \mathbb{R}^2$, one has $|\mathbb{E}[f(\xi, x)]| \leq \mathbb{E}[|g(\xi, x)|] + 35\|x\|^2 \leq |x_1| + \frac{1}{64} + 35\|x\|^2 < \infty$.

Step 3. Let $p_k = (0, \frac{1}{k}) \in \mathbb{R}^2$. Then for any $k \in \mathbb{N}$ and $y \in \mathbb{B}(p_k, \frac{r_k}{2})$, by (8), one has

$$\sum_{j=1}^{\infty} \psi_j(y_2)(2 \text{bit}_j(\xi) - 1) = \eta_k(2 \text{bit}_k(\xi) - 1), \quad \sum_{j=1}^{\infty} \psi'_j(y_2)(2 \text{bit}_j(\xi) - 1) = 0,$$

which implies that

$$g(\xi, y) = \begin{cases} y_1 + \eta_k & \text{if } \text{bit}_k(\xi) = 1, \\ y_1 - \eta_k & \text{if } \text{bit}_k(\xi) = 0, \end{cases} \quad \nabla_x g(\xi, y) = \begin{bmatrix} 1 \\ 0 \end{bmatrix}.$$

Let $\Delta_k = \frac{1}{32Ck^4} < \frac{1}{16k^2} = \frac{r_k}{2}$. For any $\xi \in \Xi, k \in \mathbb{N}$, and $y \in \mathbb{B}(p_k, \Delta_k)$, one has $|y_1| \leq \Delta_k < \frac{1}{64k^4} = \eta_k$. Thereby, it follows that

$$\text{sgn}(g(\xi, y)) = 2 \text{bit}_k(\xi) - 1, \quad \partial h(g(\xi, y)) = \{\text{bit}_k(\xi)\}, \quad (9)$$

and $f(\xi, \cdot)$ is smooth on $\cup_k \mathbb{B}(p_k, \Delta_k)$ for any $\xi \in \Xi$. Hence, for any $k \in \mathbb{N}$ and $y \in \mathbb{B}(p_k, \Delta_k)$, by [14, Exercise 10.26], we can write

$$\nabla_x f(\xi, y) = \nabla_x g(\xi, y) \partial h(g(\xi, y)) + 70y = \nabla_x g(\xi, y) \text{bit}_k(\xi) + 70y = \begin{bmatrix} \text{bit}_k(\xi) \\ 0 \end{bmatrix} + 70y,$$

where we omit curly braces for singletons for simplicity.

Step 4. Let \mathcal{F} -measurable $E \subset \Omega$ with $\mathbb{P}(E) = 1$, $\bar{\nu}$, and $\{\mathbf{k}^\nu\}_\nu$ be defined in Lemma 2. For any $\omega \in E$, $\nu \geq \bar{\nu}(\omega)$, one has $\omega \in E_{\mathbf{k}^\nu(\omega)}^\nu$ for some $\mathbf{k}^\nu(\omega) \leq \lceil 2^{\nu+1} \log(\nu+1) \rceil$. Then, $\text{bit}_{\mathbf{k}^\nu(\omega)}(\boldsymbol{\xi}^i(\omega)) = 1$ for all $i \in [\nu]$. Let $\mathbf{p}^\nu(\omega) = p_{\mathbf{k}^\nu(\omega)} \in \{p_k\}_k$. For any $y \in \mathbb{B}(\mathbf{p}^\nu(\omega), \Delta_{\mathbf{k}^\nu(\omega)})$, we have

$$\frac{1}{\nu} \sum_{i=1}^\nu \nabla_x f(\boldsymbol{\xi}^i(\omega), y) = \begin{bmatrix} \frac{1}{\nu} \sum_{i=1}^\nu \text{bit}_{\mathbf{k}^\nu(\omega)}(\boldsymbol{\xi}^i(\omega)) \\ 0 \end{bmatrix} + 70y = \begin{bmatrix} 1 \\ 0 \end{bmatrix} + 70y.$$

Meanwhile, for any $k \in \mathbb{N}$ and $y \in \mathbb{B}(p_k, \Delta_k)$, deterministically, we have

$$\mathbb{E}[\nabla_x f(\boldsymbol{\xi}, y)] = \begin{bmatrix} \mathbb{E}[\text{bit}_k(\boldsymbol{\xi})] \\ 0 \end{bmatrix} + 70y = \begin{bmatrix} 1/2 \\ 0 \end{bmatrix} + 70y.$$

Step 5. Note that $C = 70$, and

$$0 \leq \delta^\nu = \frac{1}{2240 \lceil 2^{\nu+1} \log(\nu+1) \rceil^4} = \Delta_{\lceil 2^{\nu+1} \log(\nu+1) \rceil} \leq \Delta_{\mathbf{k}^\nu(\omega)}.$$

Therefore, for any $\omega \in E$, we conclude that $\mathbf{p}^\nu(\omega) \in \{p_k \mid 1 \leq k \leq \lceil 2^{\nu+1} \log(\nu+1) \rceil\} \subset D^\nu$ and

$$\begin{aligned} & \liminf_{\nu \rightarrow \infty} \sup_{x \in D^\nu} \inf_{y, \hat{y} \in \mathbb{B}(x, \delta^\nu)} \left\| \mathbb{E}[\nabla_x f(\boldsymbol{\xi}, y)] - \frac{1}{\nu} \sum_{i=1}^\nu \nabla_x f(\boldsymbol{\xi}^i(\omega), \hat{y}) \right\| \\ & \geq \sup_{\nu' \geq \bar{\nu}(\omega)} \inf_{\nu \geq \nu'} \inf_{y, \hat{y} \in \mathbb{B}(\mathbf{p}^\nu(\omega), \delta^\nu)} \left\| \mathbb{E}[\nabla_x f(\boldsymbol{\xi}, y)] - \frac{1}{\nu} \sum_{i=1}^\nu \nabla_x f(\boldsymbol{\xi}^i(\omega), \hat{y}) \right\| \geq \frac{1}{2} - \inf_{\nu' \geq \bar{\nu}(\omega)} \sup_{\nu \geq \nu'} 140\delta^\nu = \frac{1}{2}, \end{aligned}$$

which completes the proof, since (Ω, \mathcal{F}) is \mathbb{P} -complete.

Acknowledgement. This work is supported in part by the Office of Naval Research under grant N00014-24-1-2492.

References

- [1] Z. Artstein and R. J-B Wets. Consistency of minimizers and the SLLN for stochastic programs. *Journal of Convex Analysis*, 2(1-2):1–17, 1995.
- [2] A. Brøndsted and R. T. Rockafellar. On the subdifferentiability of convex functions. *Proceedings of the American Mathematical Society*, 16(4):605–611, 1965.
- [3] F. H. Clarke. *Optimization and Nonsmooth Analysis*. SIAM, 1990.
- [4] D. Davis and D. Drusvyatskiy. Graphical convergence of subgradients in nonconvex optimization and learning. *Mathematics of Operations Research*, 47(1):209–231, 2022.
- [5] D. J. Foster, A. Sekhari, and K. Sridharan. Uniform convergence of gradients for non-convex learning and optimization. *Advances in Neural Information Processing Systems*, 31, 2018.
- [6] J.-B. Hiriart-Urruty and C. Lemaréchal. *Convex Analysis and Minimization Algorithms I: Fundamentals*, volume 305. Springer Science & Business Media, 2013.
- [7] A. D. Ioffe. An invitation to tame optimization. *SIAM Journal on Optimization*, 19(4):1894–1917, 2009.

- [8] O. Kallenberg. *Foundations of Modern Probability*. Springer, 3rd edition, 2021.
- [9] S. Mei, Y. Bai, and A. Montanari. The landscape of empirical risk for nonconvex losses. *Annals of Statistics*, 46(6A):2747–2774, 2018.
- [10] I. Molchanov. *Theory of Random Sets*. Springer, 2nd edition, 2017.
- [11] V. I. Norkin. Stochastic Lipschitz functions. *Cybernetics and Systems Analysis*, 22(2):226–233, 1986.
- [12] V. I. Norkin and R. J-B Wets. On strong graphical law of large numbers for random semicontinuous mappings. *Vestnik Sankt-Peterburgskogo Universiteta*, Seriya 10(3):102–111, 2013.
- [13] R. T. Rockafellar. *Convex Analysis*. Princeton University Press, 1970.
- [14] R. T. Rockafellar and R. J-B Wets. *Variational Analysis*. Springer, 3rd printing-2009 edition, 1998.
- [15] F. Ruan. On the uniform convergence of subdifferentials in stochastic optimization and learning. *Mathematics of Operations Research*, to appear, 2025.
- [16] A. Salim. A strong law of large numbers for random monotone operators. *Set-Valued and Variational Analysis*, 31(4):38, 2023.
- [17] A. Shapiro. Asymptotic properties of statistical estimators in stochastic programming. *Annals of Statistics*, 17(2):841–858, 1989.
- [18] A. Shapiro and Y. Wardi. Convergence analysis of stochastic algorithms. *Mathematics of Operations Research*, 21(3):615–628, 1996.
- [19] A. Shapiro and H. Xu. Uniform laws of large numbers for set-valued mappings and subdifferentials of random functions. *Journal of Mathematical Analysis and Applications*, 325(2):1390–1399, 2007.
- [20] A. Shapiro, D. Dentcheva, and A. Ruszczyński. *Lectures on Stochastic Programming: Modeling and Theory*. SIAM, 3rd edition, 2021.
- [21] P. Terán. On a uniform law of large numbers for random sets and subdifferentials of random functions. *Statistics & Probability Letters*, 78(1):42–49, 2008.
- [22] A. W. van der Vaart and J. A. Wellner. *Weak Convergence and Empirical Processes*. Springer, 2nd edition, 2023.
- [23] D. Williams. *Probability with Martingales*. Cambridge University Press, 1991.
- [24] H. Xu. Uniform exponential convergence of sample average random functions under general sampling with applications in stochastic programming. *Journal of Mathematical Analysis and Applications*, 368(2):692–710, 2010.