



A Local Nearly Linearly Convergent First-Order Method for Nonsmooth Functions with Quadratic Growth

Damek DavisLiweijiang

Received: 3 June 2022 / Revised: 15 July 2023 / Accepted: 21 November 2023 © SFoCM 2024

Abstract

Classical results show that gradient descent converges linearly to minimizers of smooth strongly convex functions. A natural question is whether there exists a locally nearly linearly convergent method for nonsmooth functions with quadratic growth. This work designs such a method for a wide class of nonsmooth and nonconvex locally Lipschitz functions, including max-of-smooth, Shapiro's decomposable class, and generic semialgebraic functions. The algorithm is parameter-free and derives from Goldstein's conceptual subgradient method.

Keywords Subgradient method \cdot Goldstein subdifferential \cdot Semialgebraic \cdot Partial smoothness \cdot VU-structure

1 Introduction

Slow sublinear convergence of first-order methods in nonsmooth optimization is often illustrated with the following simple strongly convex function:

$$f(x) = \max_{1 \le i \le m} x_i + \frac{1}{2} x^2 \quad \text{for some } m \le d \text{ and all } x \in \mathbb{R}^d.$$
 (1.1)

For example, consider the subgradient method applied to f, which generates iterates x_k . Since f is strongly convex, classical results dictate that $f(x_k)$ – inf $f = O(k^{-1})$.

Communicated by Jérôme Bolte.

Research of Davis supported by an Alfred P. Sloan research fellowship and NSF DMS award 2047637.

B Damek Davis dsd95@cornell.edu http://people.orie.cornell.edu/dsd95/ Liwei Jiang https://liwei-jiang97.github.io/

Published online: 14 June 2024

School of Operations Research and Information Engineering, Cornell University, Ithaca, NY 14850, USA

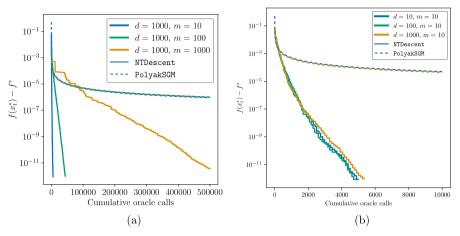


Fig. 1 Comparison of NTDescent with PolyakSGM on (1.1). Left: we fix d and vary m; Right: we fix m and vary d. For both algorithms, the value $f(x_t^*)$ denotes the best function value seen after t oracle evaluations

On the other hand, under proper initialization and an adversarial first-order oracle, there is a matching lower bound for the first m iterations: $f(x_k)$ – inf $f \ge (2m)^{-1}$ for all $k \le m$; see [11, 43]. Beyond the subgradient method, the lower bound also holds for any algorithm whose kth iterate lies within the linear span of the initial iterate and past k-1 computed subgradients. Thus, one must make more than m first-order oracle calls to f, i.e., function and subgradient evaluations, before possibly seeing improved convergence behavior.

While such methods make little progress when $k \le m$, this behavior may or may not continue for k = m. On one extreme, the subgradient method continues to converge slowly even when equipped with the popular Polyak stepsize PolyakSGM) [47]; see dashed lines in Fig. 1. On the opposite extreme, more sophisticated algorithms such as the center of gravity method or the ellipsoid method converge linearly, but their complexity scales with the dimension of the problem, a necessary consequence of the linear rate of convergence; see the discussion in [11, Chapter 2].

A natural question is whether there exists a first-order method whose behavior lies in between these two extremes, at least for nonsmooth functions f satisfying regularity conditions at local minimizers. Regularity conditions often take the form of growth—linear or quadratic—away from minimizers. Well-known results show that subgradient methods converge linearly on nonsmooth functions with linear (also called sharp) growth [47]. On the other hand, in smooth convex optimization, quadratic growth entails linear convergence of gradient methods. However, to the best of our knowledge, no parallel result for nonsmooth functions with quadratic growth exists. Thus, in this work, we ask

is there a locally nearly linearly convergent method for nonsmooth functions with quadratic growth whose rate of convergence and region of rapid local convergence solely depends on *f* ?

Let us explain the qualifiers "nearly" and "solely depends on f." First, the qualifier "nearly" signifies that the method locally achieves a function gap of size \mathcal{E} using at most, say, $O(C_f \log^3(1/\mathcal{E}))$ first-order oracle evaluations of f, where C_f depends on f. Second, the qualifier "solely depends on the function," signifies that f and the size of the region of local convergence do not depend on the dimension of the problem, but instead depend only on the function f through intrinsic quantities, such as Lipschitz and quadratic growth constants.

In this work, we positively answer the above question for a class of nonsmooth optimization problems with quadratic growth. The method we develop is called *Normal Tangent Descent* (NTDescent). We formally describeNTDescent in Sect. 1.4. For now, we illustrate the performance of NTDescent on f from (1.1) in Fig. 1. In both plots, we see NTDescent improves on the performance of PolyakSGM, measured in terms of oracle calls. This is a fair basis for comparison since botholyakSGM and NTDescent perform a similar amount of computation per oracle call. Figure 1b also shows that the performance of NTDescent is dimension independent. We highlight that this performance was achieved without any tuning of parameters for TDescent. Indeed, our main theoretical guarantees for NTDescent (Theorem 1.1) do not require the user to set any parameters.

The problem class on which NTDescent succeeds consists of locally Lipschitz nonsmooth functions with quadratic growth and a certain smooth substructure at local minimizers. Importantly, we do not assume the problems under consideration are convex, though convexity entails improved guarantees. Two example classes with such smooth substructure include (i) "generic" semialgebraic functions and (ii) properly C^p decomposable loss functions satisfying strict complementarity and quadratic growth conditions [49]. A semialgebraic function is one whose graph is the finite union of intersections of polynomial inequalities. Semialgebraic functions (more generally tame [32] functions) model most problems of interest in applications. If f is semialgebraic, for a full Lebesgue measure set of $W \in \mathbb{R}^d$, we will use show that the tilted function $f w : x \to f(x) + W$ x has quadratic growth and the desired smooth substructure at each local minimizer, explaining the qualifier "generic." We mention that this fact essentially follows from combining results of [19, 23]. On the other hand, a properly C^p decomposable function is one that decomposes near local minimizers as a composition of a positively homogeneous convex function with a smooth mapping that maps the minimizer to the origin. Decomposable functions appear often in practice, e.g., in eigenvalue and data fitting problems. An important subclass of decomposable functions consists of so-called "max-of-smooth" functions, which are the maximum of finitely many smooth functions that satisfy certain regularity conditions at minimizers, e.g., f in (1.1).

The precise smooth substructure used in this work was recently identified in [19], where it was shown to be available in decomposable and generic semialgebraic problems. Since it is available in many problems of interest, throughout this introduction we call this the combination of quadratic growth and smooth substructure *typical structure* and call functions possessing this combined structure *typical*. We present the formal structure in Sect. 3. At the heart of this structure is a distinguished smooth manifold M—called the *active manifold*—containing a local minimizer of interest. We formally define the active manifold concept in Definition 1.2, but at a high level, the

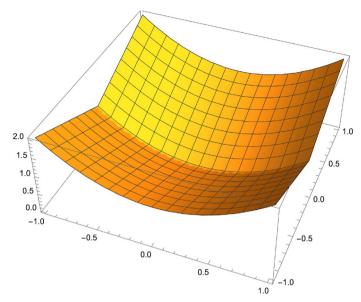


Fig. 2 The function $f(u, v) = u^2 + |v|$ has typical structure

two crucial characteristics are that (i) along the manifold, the function fis smooth and (ii) normal to the manifold, the function grows sharply. For example, Fig. 2 depicts the nonsmooth function $f(u, v) = u^2 + |v|$ for which the *u*-axis plays the role of *M* . In Sect. 1.3.2 we will examine this function and explain how we use its typical structure in NTDescent. This example also has the smooth substructure developed in several seminal works in the optimization literature, including those found in work on identifiable surfaces [52], partly smooth functions [38], VU-structures [36, 41], and minimal identifiable sets [26]. However, crucial to the analysis of NTDescent are two further properties introduced in [19], called *strong* (a)-regularity and $(b \le)$ *regularity*. Strong (a)-regularity roughly states that the function is smooth in tangent directions to the manifold up to an error term which is linear in the distance to the manifold. On the other hand, (b_{\leq}) -regularity is a one-sided uniform semismoothness [40] property that holds automatically when fis (weakly) convex. Both properties hold for the two variable example in Fig. 2 and for the function in (1.1), where the active manifold is the subspace in which the first *m* variables take on the same value: $M = \{ x \in \mathbb{R}^d : x_1 = x_2 = \dots x_m \}.$

Before turning to the description of NTDescent, we point out that similar smooth substructure has been used in the analysis first-order methods in nonsmooth optimization, most famously for functions with VU-structure [36, 41] and more recently for max-of-smooth functions. For VU functions, so-called "bundle-methods," [35, 50] which possess an inner-outer loop structure, have been shown to converge superlinearly with respect to the number of outer-loop steps [41]; see also the survey [44]. These methods have excellent empirical performance, but a complete account of their

¹ Though they also benefit from smooth substructure, *proximal-methods* do not fall within the oracle model of first-order methods considered in this work. Thus, we omit them from our discussion.

inner-loop complexity remains elusive. On the other hand, in a recent work, Han and Lewis proposed a first-order method—Survey Descent—that converges linearly on certain strongly convex max-of-smooth objectives, stepping beyond the classical smooth setting [30]. The method shows favorable performance beyond the max-of-smooth class, e.g., on certain eigenvalue optimization problems, but no theoretical justification for this success is available. We discuss Survey Descent in more detail in Sect. 7.1. We now motivateNTDescent.

1.1 Motivation: Goldstein's Conceptual Subgradient Method

To motivate NTDescent and the role of smooth substructure, let us set the stage: consider the nonsmooth optimization problem:

minimize_{$$x \in \mathbb{R}^d$$} $f(x)$,

where $f: \mathbb{R}^d \to \mathbb{R}$ is a locally Lipschitz function, which is not necessarily convex. The algorithm developed in this work assumes *first-order oracle access* to f [11, 42, 43]. In particular, at every $x \in \mathbb{R}^d$ we must be able to evaluate f(x) and retrieve an element of the *Clarke subdifferential* $\partial f(x)$. Informally, the Clarke subdifferential is comprised of convex combinations of limits of gradients taken at nearby points; a formal definition appears in Sect. 1.7. The Clarke subdifferential reduces to the familiar objects in classical settings. For example, when f is C^1 , the Clarke subdifferential reduces to the singleton mapping $\{\nabla f\}$. In addition, when f is convex, the Clarke subdifferential reduces to the subdifferential in the sense of convex analysis.

The starting point of this work is the classical conceptual subgradient method of Goldstein [29]. The core object in this method is the Goldstein subdifferential:

$$\partial_{\sigma} f(x) := \operatorname{conv} \left(\begin{array}{c} f(y) \\ y \in \overline{B}\sigma(x) \end{array} \right) \quad \text{for all } x \in \mathbb{R}^d \text{ and } \sigma > 0.$$
 (1.2)

This subdifferential is simply the convex hull of all Clarke subgradients of f taken at points inside the ball of radius σ . Its importance arises from the following descent property proved in [29]: fix $\sigma > 0$ and $x \in \mathbb{R}^d$ and let W denote the minimal norm element of $\partial_{\sigma} f(x)$. Then

$$f \quad x - \sigma \frac{w}{w} \leq f(x) - \sigma w \qquad \text{if } w = 0. \tag{1.3}$$

This property motivates Goldstein's conceptual subgradient method, which simply iterates:

$$x_{k+1} = x_k - \sigma \frac{w_k}{w_k}$$
 where $w_k = \underset{w \in \partial \sigma f(x_k)}{\operatorname{argmin}} w.$ (1.4)

This algorithm is remarkable since it is provably a descent method for any Lipschitz function and even converges at a sublinear rate. Indeed, a quick appeal to (1.3) yields

$$\min_{k=0,\dots,K-1} w_k \le \varepsilon \qquad \text{holds when} \qquad K \ge \frac{f(x_0) - \min f}{\sigma \varepsilon}.$$

While this exact variant of the Goldstein method is not necessarily implementable, recent work has devised approximate versions of the method that have similar sublinear convergence properties [21, 55].

The algorithm introduced in this work approximately implements the method (1.4). The goal of this work is to prove that the method is locally nearly linearly convergent on typical nonsmooth functions. To develop such a method, we must resolve two issues for this problem class. First, we must develop rapidly convergent algorithms that approximately compute the minimal norm element of the Goldstein subdifferential. Second, we must devise an appropriate regularity property that ensures the proposed method converges nearly linearly. We will discuss both of these properties in turn, beginning with a regularity property that relates the decrement in (1.3) to the function gap.

1.2 Linear Convergence via a Gradient Inequality

Observe that if the bound

$$\sigma w_k \ge \eta (f(x_k) - \min f)$$

holds for some $\eta > 0$ and all k > 0, then the Goldstein method (1.4) converges linearly to a minimizer of f. A potential issue with this inequality is that the vector W_k is zero whenever σ is larger than the distance of x_k to the nearest critical point of f; thus the algorithm may stall whenever x_k is near enough to a minimizer. This suggests a simple relaxation of the property that allows σ to depend on x_k .

Indeed, we will provide conditions under which the following bound holds near a local minimizer \bar{x} of f: there exists a constant $\eta > 0$ and a function $\sigma : \mathbb{R}^d \to \mathbb{R}_+$ such that for all x near \bar{x} , we have

$$\sigma(x)\operatorname{dist}(0,\,\partial_{\sigma(x)}f(x)) \ge \eta(f(x) - f(x)). \tag{1.5}$$

throughout, we will refer to this bound as a *gradient inequality*, due to its similarity to the Kurdyka-Łojasiewicz (KL) gradient inequality [7]. The KL inequality requires that a suitable nonlinear reparameterization $\psi: \mathbb{R} \to \mathbb{R}$ of the function gap is bounded by the minimal norm Clarke subgradient for all x near x:

$$\operatorname{dist}(0,\,\partial\,f(x))\geq\psi(\ f(x)-\ f(\vec{x})).$$

In recent years, the KL inequality has played a key role in establishing convergence and rates of convergence for proximal methods in nonsmooth optimization and in continuous time analogs of the subgradient method; see e.g., [3, 4, 7, 9, 53].

To illustrate, let us specialize to the semialgebraic setting, where the desingularization function Ψ is known to take the form $\Psi(r) = r^{\theta}$ for $\theta \in [0, 1)$. The work [2, Theorem 2] initiated the study of convergence of proximal methods in this setting, showing that the proximal point method asymptotically converges to its limit point, which is critical but not necessarily optimal. The method convergence in finitely many steps when $\theta = 0$, locally converges linearly when $\theta \in (0, 1/2]$, and locally converges at the rate $k^{\frac{-(1-\theta)}{2\theta-1}}$ when $\theta \in (1/2, 1)$. Further works such as [3, 9] generalized the techniques to related proximal methods. Passing to continuous time, one is interested in the convergence of the trajectory of subdifferential inclusion satisfying $\dot{x}(t) \in -\partial f(x(t))$ at almost every *t* . Here, the rates of convergence exactly parallel those in the proximal methods as shown in [6, Theorem 4.7]. In contrast to the proximal and continuoustime settings, we do not know whether the KL inequality alone allows one to design a locally linearly convergent discrete-time subgradient method, except in the setting where $\theta = 0$ (i.e., f is sharp) and f is convex [47] or weakly convex [22]; weakly convex functions form a broad class of nonconvex functions that includes all compositions of Lipschitz convex functions with smooth mappings. When $\theta \geq 0$, to the best of our knowledge, the best rate proved in the literature for any subgradient type method is $k^{\frac{-(1-\theta)}{2\theta}}$ [33]; this result is only known to hold for convex functions.³

A well-known property of the KL inequality is its prevalence: it holds at each critical point of an arbitrary lower-semicontinuous semialgebraic function f [7]. We will show that the gradient inequality (1.5) is also prevalent in the sense that it holds for the aforementioned problems with typical structure. In this way, the conceptual method (1.4) with varying $\sigma_k := \sigma(x_k)$ will locally converge linearly on such problems. The reader may wonder whether we can or must find the precise value $\sigma(x_k)$. We will show that for typical problems, an appropriat $\sigma(x_k)$ may be found through a line search procedure.

1.3 Approximately Implementing Goldstein's Method

The gradient inequality ensures that the conceptual Goldstein method converges linearly, provided the stepsize σ is chosen adaptively. To move beyond the conceptual setting, we must develop strategies for approximating the minimal norm element of $\partial_{\sigma} f(x)$ for $\sigma > 0$ and $x \in \mathbb{R}^d$. Let us suppose we have such a method and denote it by MinNorm(x, σ). Then the method of this work simply iterates:

$$x_{k+1} = x_k - \sigma_k \frac{w_k}{w_k}$$
 and $w_k = \text{MinNorm}(x_k, \sigma_k)$ (1.6)

for an appropriate sequence $\sigma_k > 0$. We will discuss and develop two different implementations of MinNorm(x, σ) in this work. Given $x \in \mathbb{R}^{d}$ and $\sigma > 0$,

 $^{^2}$ These rates were shown only for "lower- C^2 " semialgebraic losses, but extend to locally Lipschitz semialgebraic functions via the semialgebraic "chain rule" proved in [20].

³ The results stated in [33] pertain to functions with Hölder growth; thus, to prove the results stated in the paragraph, we must use the following known fact: functions satisfying the KL inequality with exponent θ have Hölder growth with exponent $1/(1-\theta)$, which follows from the proof of [24, Theorem 3.7].

both methods iteratively construct a sequence of Clarke subgradients g_0, \dots, g_{T-1} taken at points in the ball $\overline{B}\sigma(x)$ and then output a "small" convex combination $W \in \text{conv}\{g_0, \dots, g_{T-1}\}$, which satisfies the descent condition

$$f \quad x - \sigma \frac{w}{w} \le f(x) - \frac{\sigma}{8} w. \tag{1.7}$$

The oracle complexity of $MinNorm(x, \sigma)$ is then T function/subgradient evaluations, and we hope to ensure that T is relatively small, for example, a constant or at most

$$T = O \log_{x,\sigma}^{-1}$$
 where $_{x,\sigma} := \operatorname{dist}(0, \partial_{\sigma} f(x)).$

Provided that T is on this order, that f satisfies the gradient inequality (1.5), and that σ_k is chosen appropriately, the iterate x_k will satisfy $f(x_k) - f(\overline{x}) \le \varepsilon$ after at most $O(\log^2(1/\varepsilon))$ iterations, a nearly linear rate of convergence. This complexity ignores the cost of choosing an appropriate stepsiz \mathcal{O}_k , but we will show that in typical problems we can find appropriate σ_k with at most $O(\log(1/\varepsilon))$ function/subgradient evaluations.

We are aware of two MinNorm type methods in the literature, but their complexity is either too large or is useful only in low dimensions problems. For example, the works [21, 55] introduced such a method for general locally Lipschitz functions. However, the complexity of the method is $T = O(1/x, \sigma)$ —too large for our purposes. On the other hand, the work [21] also introduced a method tailored to low-dimensional weakly convex functions. However, the method is based on cutting plane techniques, so its complexity scales linearly with dimension: $T = O(d \log(1/x, \sigma))$.

While existing MinNorm methods are slow for general Lipschitz functions, we show that the aforementioned typical structure allows us to develop MinNorm methods that accelerate in a neighborhood of the minimizer. Our approach is based on a decomposition of a neighborhood of the minimizer into two regions: one where the method of [21, 55] is applicable, and another region where a noveMinNorm method may be applied.

1.3.1 The Normal and Tangent Regions

In this work, we use the active manifold M to split the space of (x, σ) for x nearby the minimizer $x \in M$ into two sets where fast MinNorm methods are available. We call the first set the *normal region*. This region consists of points whose normal distance dist(x, M) is larger than a multiple of the squared tangential distance $PM(x) = x^{-2}$, together with stepsizes σ proportional to a multiple of the normal distance:

$$\begin{array}{ll} \frac{a_1}{2} \mathrm{dist}(x,\,M) \leq \sigma \leq & a_1 \mathrm{dist}(x,\,M); \\ a_2^2 & PM(x) - \overline{x}^2 \leq \mathrm{dist}(x,\,M), \end{array}$$

for problem dependent constants a_1 , $a_2 \in (0, 1)$; see Theorem 4.3 for more details. We will show that in this region, we have $a_1, a_2 = a_1$, so the MinNorm method of [21, 55] terminates with descent in finitely many steps.

On the other hand, we call the second set the *tangent region*. This set consists of points whose squared tangential distance is larger than a multiple of the normal distance, together with stepsizes proportional to a multiple of the tangential distance:

$$\frac{a_2}{2} P_M(x) - \overline{x} \le \sigma \le a_2 P_M(x) - \overline{x};$$

$$\frac{\operatorname{dist}(x, M)}{\sigma} \le 2a_2 P_M(x) - \overline{x},$$

where a_1 and a_2 are as in the normal region. For this region, we will propose a new MinNorm method, which terminates rapidly. We note that in both cases we provide a range of valid σ , rather than a single value since we aim to estimate σ with a line search.

1.3.2 A Simple Example

Before describing the methods in detail, let us illustrate the regions and the principles of the methods on the following simple function of two variables $f(u, v) = u^2 + |v|$, which has a unique minimizer $a\bar{x} = (0, 0)$. Here, the u-axis is the active manifold M. Along the manifold, f is smooth and grows quadratically, while off of the manifold f grows sharply; see Fig. 2 for a plot of the function and see Fig. 3 for the x-component of the normal and tangent regions for f (with $a_2 = 1/8$). The manifold f induces a decomposition of f into smooth f f f f into smooth

$$f(x) = fU(x) + fV(x) = PM(x) - x^{2} + dist(x, M).$$
 (1.8)

From this decomposition, we see fv is dominant in the normal region, while fu is dominant in the tangent region. Likewise, as we will argue momentarily, the minimal norm Goldstein subgradient $W_{\sigma} \in \partial_{\sigma} f(x)$ satisfies $w_{\sigma} \ge \nabla fv(x)$ in the normal region, while $w_{\sigma} = (\nabla fu(x))$ in the tangent region. This has several consequences. First, in the normal region, the MinNorm method of [21, 55] will terminate in finitely many steps, due to the lower bound $w_{\sigma} \ge 1$. On the other hand, in the tangent region, w_{σ} can be much smaller, so we must introduce a new method to generate descent. Finally, assuming these approximations are accurate, the gradient inequality (1.6) quickly follows: in the normal region, we have

$$\sigma w \sigma = di(st(x, M)) = (f(x)),$$

while in the tangent region, we have

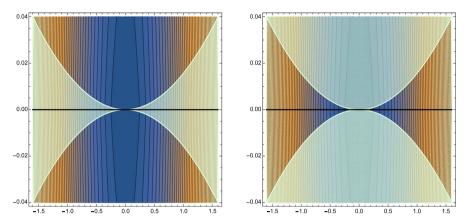


Fig. 3 Contour plots for $f(u, v) = u^2 + |v|$ together with M, shown in black. The x-component of the tangent (left) and normal (right) regions are overlaid in light green

Though it follows from immediate calculations in this example, in the more general setting, the following consequence of quadratic growth will be crucial in establishing a similar bound: $\nabla f_U(x) = (P_M(x) - \overline{x})$.

Now, to lower bound $W \sigma$ we use the following fact: $\nabla f U(u, V)$ is tangent to M, while $\nabla f V(u, V)$ is normal to M when V = 0. Thus, to lower bound $W \sigma$ in the normal region, we will simply lower bound the size of the normal component of W_{σ} . Indeed, since $\sigma < \operatorname{dist}(x, M)$, all points $x \in B_{\sigma}(x)$ are on the same side of M. Therefore, the normal component of W_{σ} is an average of *identical* gradients $\nabla f V(x) = \nabla f V(x)$. Likewise, in the tangent region, we lower bound the tangent component of W_{σ} . Indeed, since $\sigma < \overline{x} - P_M(x)/8$, the projection onto M of all points $x \in B_{\sigma}(x)$ are on the same side of the origin. Thus, the tangent component of W_{σ} is an average of nearly identical gradients $\nabla f U(x) \approx \nabla f U(x)$, yielding the lower bound. We prove a more general form of both of these lower bounds Lemma 4.1 and Lemma 4.2, which follow from similar principles.

Turning to algorithms, we have so far noted that the MinNorm method of [21, 55] may be used in the normal region. In the tangent region, we are unsure how to design a method that can quickly recover $W\sigma$. Instead of searching for $W\sigma$ directly, we take a slightly different perspective in the tangent region: we seek a vector $g \in \partial \sigma f(x)$ with "small" normal component, meaning:

$$P_N(g) = O(\nabla f \cup (x)^2)$$

where N is the normal space to M. Intuitively, when g has small normal component, the nonsmooth part fv minimally changes along a gradient step. On the other hand, if g is sufficiently correlated with ∇fu , the smooth part fu decreases at an appropriate rate; we prove this in a more general setting in Lemma 5.2.

Why might one expect such a *g* to be available in the tangent region? The reason is that the gradient of the smooth component is itself a Goldstein subgradient. Indeed, for points near the origin and in the tangent region, the tangential distance is much larger

than the normal distance. Thus, the reflection of any point V across the manifold is contained in $B\sigma(x)$, which immediately implies gradient of the smooth component is an element of Goldstein subdifferential:

$$\nabla \ f \cup (u, \ v) = \frac{1}{2} \nabla \ f(u, \ v) + \frac{1}{2} \nabla \ f(u, \ -v) \in \partial \ \sigma \ f(x). \tag{1.9}$$

While the inclusion (1.9) illustrates one way to construct such a g, we cannot hope for perfect symmetry in general problems.

Instead, a central insight of this work is that a similar approximate reflection exists in problems with typical structure. To illustrate, consider Fig. 4. This figure depicts a point x in the tangent region together with the result of a normalized gradient step:

$$x_+ := x - \sigma \frac{\nabla f(x)}{\nabla f(x)}$$
.

As can be seen from the figure, x_+ is an approximate reflection of x across the u-axis, which "flips the sign" of the nonsmooth component of $\nabla f: \nabla f_V(x) = -\nabla f_V(x_+)$. Thus, in this setting, one may "cancel out" the nonsmooth component by a simple averaging:

$$\nabla \ f \upsilon(x) \approx \ \frac{1}{2} \nabla \ f(x) + \ \frac{1}{2} \nabla \ f(x_+).$$

While seemingly crude, we will show this strategy generalizes to typical functions. An important distinction with the general setting is that a single averaging step alone will no longer suffice. Nevertheless, we show that by iterating this process, we can geometrically shrink the normal component of the Goldstein gradient, eventually yielding descent.

1.3.3 Two MinNorm Methods: NDescent and TDescent

To generalize the strategy outlined in the previous section, we will prove that the minimal norm Goldstein subgradients of typical problems similarly split into tangent and normal components just as in Sect. 1.3.2. Then, we introduce twoMinNorm type methods for "normal" and "tangent" steps.

For (x, σ) in the normal region, we use a small modification of the MinNorm type method of [21]. We call this method *Normal Descent* (NDescent) and describe it in Algorithm 1. As in the simple example above, we will show that NDescent must terminate with an approximately minimal norm Goldstein subgradient in finitely many steps, provided σ lies within an appropriate range. We will show that this subgradient is a descent direction satisfying (1.7).

Algorithm 1 NDescent (x, g, σ, T)

```
1: Set g_0 = g and t = 0.
```

2: while $T-1 \ge t$, $g_t > 0$, and $\frac{\sigma}{8} g_t \ge f(x) - f(x - \sigma \frac{g_t}{q_t})$ do

3: Choose any *r* satisfying $0 \le r \le \sigma$ g_t .

4: Sample ζ_t uniformly from $B_r(g_t)$.

5: Choose y_t uniformly at random in the segment x_t $x - \sigma \frac{\zeta_t}{\zeta_t}$.

6: Choose $\hat{g}_t \in \partial f(y_t)$.

7: $g_{t+1} = \operatorname{argmin}_{z \in [g_t, \hat{g}_t]} z_2$.

8: t = t + 1.

9: end while

10: return g_t .

We illustrate the principle behind NDescent as follows. Suppose we are given a vector $g \in \partial_{\sigma} f(x)$ not satisfying the descent condition, i.e., with $u := \frac{g}{a}$, we have

$$f\left(x-\sigma u\right)-f\left(x\right)\geq-\frac{g}{8}.$$

Then by Lebourg mean value theorem [16, Theorem 2.4] (provided that f is differentiable along the line segment between [x, x], which can be ensured by adding a small perturbation to g; we ignore this in our discussion), we may assume that

$$f(x - \sigma u) - f(x) = \sigma \int_{0}^{1} - \nabla f(x - \sigma t u), u dt = -\sigma \lor, u,$$

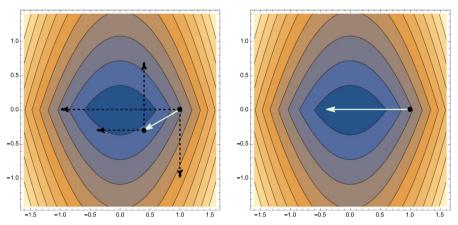


Fig. 4 Contour plots for $f(u, v) = u^2 + |v|$. Left: The point $x = (1, \cdot 1)$ together with the approximate reflection $x_+ = x - .3 \frac{\nabla f(x)}{\nabla f(x)}$ across the u axis. The solid light green arrow is parallel to the negative gradient direction $-\nabla f(x)$. The dashed arrows denote the orthogonal decomposition of $-\nabla f(x)$, respectively $-\nabla f(x_+)$, into the vectors $-\nabla f_U(x)$ and $-\nabla f_V(x)$, respectively $-\nabla f_U(x_+)$ and $-\nabla f_V(x_+)$. From the plot, we see $\nabla f_V(x) = -\nabla f_V(x_+)$. Right: The point x with estimate $-\frac{1}{2}(\nabla f(x) + \nabla f(x_+))$ of the vector $-\nabla f_U(x)$

where $v := \int_0^1 \nabla f(x - t u) dt \in \partial \sigma f(x)$. Consequently, $V, g \le g^{-2}/8$. While it is not possible to compute V, we can compute a *random* element of the Goldstein subdifferential, satisfying the same inequality in expectation. Indeed, defining $V = \nabla f(y)$ where y is uniformly sampled from the line segment $[x, x - \sigma u]$ (with end points x and $x - \sigma u$), we have $E_y[v]$, $g \le g^{-2}/8$. Based on this bound, a quick calculation shows that the minimal norm element g of the line segment [g, V] satisfies the bound

$$E_y g_+^2 \le g^2 - \frac{g^4}{16L^2}$$
.

Moreover $g_+ \in \partial_\sigma f(x)$. Thus, repeating this process yields a decreasing sequence of Goldstein subgradients which tend to zero as long as the descent condition is not met. In general, the norms of the subgradients generated by this process decay at a rate of 1/k. However, we will prove that $\operatorname{dist}(0, \partial_\sigma f(x))$ is bounded below by a fixed constant when (x, σ) is in the normal region described in Sect. 1.3.1. Consequently, the loop must exist in finite time with descent (with high probability), for otherwise we will have found a subgradient norm strictly smaller than dist $(0, \partial_\sigma f(x))$; see Proposition 5.1. The reader interested in the formal calculations may consult [21, 55].

On the other hand, for (x, σ) in the tangent region, we develop a new MinNorm type method, which likewise relies on an approximate reflection property. We call this method *Tangent Descent* (TDescent) and present it in Algorithm 2. Given an input point x, stepsize $\sigma > 0$, and initial subgradient $g_0 \in \partial f(x)$, TDescent repeats the following steps

Choose:
$$\hat{g}_k \in \partial f \quad x - \sigma \frac{g_k}{g_k}$$
;

Update:
$$g_{k+1} = \underset{g \in [g_k, \hat{g}_k]}{\operatorname{argmin}} g$$
,

until it achieves descent $f(x - \sigma \frac{g_k}{g_k}) \le f(x) - \frac{\sigma}{8} g_k$ or runs over budget.

Algorithm 2 TDescent(x, g, σ, T)

```
1: Set g_0 = g and t = 0.

2: while T - 1 \ge t, g_t > 0, and \frac{\sigma}{8} g_{t-2} \ge f(x) - f(x - \sigma \frac{g_t}{g_t}) do

3: Choose \hat{g}_t \in \partial f(x - \sigma \frac{g_t}{g_t}).

4: g_{t+1} = \operatorname{argmin}_{z \in [g_t, \hat{g}_t]} z.

5: t = t + 1.

6: end while 7: return g_t.
```

The motivation for this method is that for typical problems the step $x - \sigma \frac{g_k}{g_k}$ is locally an approximate reflection across M that "flips" the normal component of the Goldstein subgradient. Indeed, let y := PM(x) denote the projection of x onto M

and let N := NM (y) denote the normal space to M at y. Then we will prove that for all k, we have

$$P_N g_k, \hat{g}_k \leq -C P_N g_k + O(y - \bar{x}^2),$$

for some $C \ge 0$, provided σ lies within an appropriate range. This inequality ensures that each step of the TDescent geometrically decreases the "normal component" of g_k , until we arrive at a Goldstein subgradient with normal component on the order of $O(y - \overline{x}^2)$; see Sect. 5.2.2. Moreover, given $g \in \partial \sigma f(x)$ satisfying

$$P_N(q) \leq C_3 y - \overline{x}^2$$

for a particular problem dependent constant $C_3 > 0$, we will prove the descent condition

$$f x - \sigma \frac{g}{g} \le f(x) - \frac{\sigma g}{8}$$

holds; see Lemma 5.2. Combining these two facts shows that TDescent will rapidly terminate with descent.

1.4 The NTDescent Algorithm

We call the main algorithm of this work *Normal Tangent Descent* NTDescent) and present it in Algorithm 4. At a high level the method is an approximate implementation of Goldstein's conceptual subgradient method as in (1.6), using NDescent and TDescent as MinNorm type methods. As input it takes three parameters: an initial point x; a sequence of grid-sizes $\{G_k\}$ for the line search on σ ; and a sequence of budgets $\{T_k\}$ for the MinNorm type methods NDescent and TDescent. Later we will show that the user may simply set $T_k = G_k = k + 1$ for all $k \ge 0$.

Algorithm 3 linesearch(x, g, s, G, T)

```
1: Set V_0 = g.

2: for i = 0, \dots, G - 1 do

3: \sigma_i = 2^{-(G-i)}.

4: u_i = \text{TDescent}(x, v_i, \sigma_i, T).

5: V_{i+1} = \text{NDescent}(x, u_i, \sigma_i, T).

6: end for

7: \tilde{x} := \operatorname{argmin}\{f(x) : x \in \{x\} \cup \{x - \sigma_i \frac{V_{i+1}}{V_{i+1}} : \sigma_i \leq \frac{V_{i+1}}{S}, i = 0, \dots, G - 1\}\}.

8: return \tilde{x}.
```

Algorithm 4 NTDescent $(x, g, c_0, \{G_k\}, \{T_k\})$

```
Require: sg = 0, c_0 \in (0, 1]

1: Set x_0 = x and g_0 = g.

2: for k = 0, 1, \dots do

3: x_{k+1} = \text{linesearch}(x_k, g_k, \max\{g_k, c_0 g_0\}, G_k, T_k).

4: Choose g_{k+1} \in \partial f(x_{k+1}).

5: end for
```

The workhorse of NTDescent is the line search procedure in Algorithm 3 (linesearch). Let us briefly comment on the structure of this method. through 6 of Algorithm 3 implement a line search on σ . Line 7 chooses the Goldstein subgradient that provides the most descent while enforcing the trust-region constraint $\sigma_i \leq \frac{V_{i+1}}{c}$. Line 7 also ensures the NTDescent is a descent method. Within the line search procedure, we evaluate TDescent and NDescent a total of G times each. Not all of the calls to TDescent and NDescent will succeed with descent within the allotted budget T, but we will show that for typical problems, at least one will generate sufficient descent provided x_k is close enough to a local minimizer and *T* is sufficiently large. The reason at least one will succeed with descent is that given any *x* sufficiently near the solution and parameters *G* and *T* sufficiently large, linesearch will find a σ such that (x, σ) is in either the normal or tangent region described in Sect. 1.3.1. The line search allows the possibility that θ is as large as 1/2, which might force x_{k+1} to leave the region surrounding the minimizer \bar{x} . This concern is what motivates the somewhat unusual structure of the line search method wherein the MinNorm-type methods are nested. Indeed, on the one hand, the nesting ensures the norms of the Goldstein subgradients v_{i+1} are decaying as σ_i increases. On the other hand, the trust region constraint ensures that σ_i is not chosen too large, which we need for two technical reasons in our analysis: (i) it prevents x_{k+1} from leaving a small neighborhood around the minimizer where our regularity assumptions hold; (ii) we can only ensure TDescent terminates quickly when $\sigma \leq \delta_{Grid}$, for a certain radius δ_{Grid} defined in Lemma 5.7, which may be substantially smaller than 1/2.

Computationally, it may at first seem desirable to drop the trust region constraint. Figure 5 shows this may not be the case. We suspect the reason is two-fold: First, the trust region constraint allows us to cut off a range of σ from our search, which might otherwise waste oracle calls; indeed, since v_{i+1} is nonincreasing in i, and σ_i is increasing, once the trust region is violated, it will be violated for all larger i. Second, although we may take longer steps by disabling the trust region constraint, the amount of descent we expect is on the order of $v_i \in \mathcal{P}$. Thus, since the norms v_{i+1} are nonincreasing, larger stepsizes σ_i do not necessarily translate to larger descent.

Finally, we comment on our motivation for choosing the scaling $s_k = \max\{g_k, c_0 g_0\}$ in the trust region constraint. First, note that it is possible to prove, using identical techniques, that the NTDescent converges when one replaces s_k by any positive sequence that is bounded from above and below by positive constants. For our particular choice of s_k , the term $c_0 g_0$ ensures the sequence is bounded below, while the local Lipschitz continuity of f ensures that s_k is bounded above. Second, we wish for the trust region constraint to be unaffected by rescalings of f. Our particular choice of s_k guarantees scaling invariance, since the subgradients

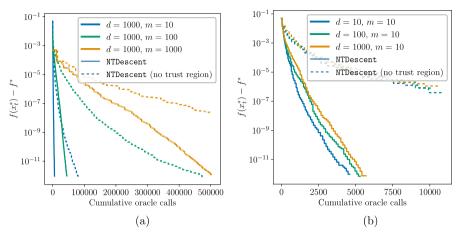


Fig. 5 Comparison of NTDescent on Problem (1.1) with the trust region constraint in Line 6 of Algorithm 3 removed. Left: we fix d and vary m; Right: we fix m and vary d. We invite the reader to compare these plots with Fig. 1

of a f are simply the subgradients of f scaled by a for any positive constant a. One might introduce other schemes for choosing s, but we did not explore such strategies. Finally, we found that performance of NTDescent is relatively insensitive to the choice of $c_0 > 0$, and any $c_0 \in \{10^{-i} : i = 0, 2, 4, 6\}$ yielded adequate performance; see Fig. 6e, f.

1.5 Main Convergence Guarantees for NTDescent

Theorem 1.1 (Main convergence theorem) Let $f: \mathbb{R}^d \to \mathbb{R}$ satisfy Assumption A at a local minimizer $x \in \mathbb{R}^d$. Fix scalar $c_0 \in (0, 1]$, budget $\{T_k\}$ and grid size $\{G_k\}$ sequences satisfying

$$\min\{T_k, G_k\} \ge k + 1$$
 for all $k \ge 0$.

Suppose that for initial point $x_0 \in \mathbb{R}^d$, there exists a subgradient $g_0 \in \partial f(x_0)$ such that $g_0 = 0$. Consider iterates $\{x_k\}$ generated by NTDescent $\{x_0, g_0, c_0, \{G_k\}, \{T_k\}\}$.

For any q, k_0 , C > 0, let $E_{k_0,q,C}$ denote the event:

$$f(x_k) - f(\bar{x}) \le \max\{(f(x_{k_0}) - f(\bar{x}))q^{k-k_0}, Cq^k\} \text{ for all } k \ge k_0.$$

Then there exists $q \in (0, 1)$, C, $C \ge 0$, and a neighborhood U of \overline{x} depending solely on f such that for any failure probability $p \in (0, 1)$ and all $k_0 \ge C \max\{\log(1/p), 1\}$, we have

$$P(E_{k_0,q,C} \mid x_{k_0} \in U) \ge 1 - p$$

provided $P(x_{k_0} \in U) > 0$. Moreover, if f is convex, we have

$$P(E_{k_0,q,C}) \ge 1 - p.$$

The theorem, which is justified in Theorems 6.3 and 6.5, bounds the function gap and distance by a quantity that geometrically decays in k. Let us examine the local complexity. Recall that each outer iteration of NTDescent requires at most $2T_kG_k$ first-order oracle evaluations. Thus, if $T_k = G_k = k + 1$ for all $k \ge 0$, the total number of oracle evaluations of K steps of NTDescent is at most $O(K^3)$. In other words, the local complexity of achieving an ε optimal solution is $O(\log^3(1/\varepsilon))$ for all sufficiently small $\varepsilon > 0$, where the big-O notation hides terms depending on the local conditioning of f; see Lemma 6.6. Therefore theorem establishes a local nearly linear rate of convergence for NTDescent.

1.6 Outline

The outline of this paper is as follows. In Sect. 1.7 we present notation and basic constructions. This section describes a key structure—the active manifold—and cannot be skipped. In Sect. 2, we present the sublinear convergence guarantees, which will be useful in the convex setting. This section also introduces key properties of the NDescent method, which will be used later in the work. In Sect. 3, we introduce our main structural assumption—Assumption A – and show that it is satisfied for the generic semialgebraic and decomposable problem classes. In Sect. 4, we show that Assumption A implies the gradient inequality (1.5). In Sect. 5 we show that the TDescent and NDescent methods terminate rapidly under appropriate conditions. In Sect. 6, we use the gradient inequality (1.5) and Assumption A to prove that NTDescent locally nearly linearly converges. Finally, in Sect. 7 we provide a brief numerical illustration.

1.7 Notation and Basic Constructions

We use standard convex analysis notation as set out in the monographs [16, 48]. Throughout, \mathbb{R}^d denotes a d-dimensional Euclidean space with the inner product \cdot , \cdot and the induced norm $x = \overline{x, x}$. We denote the open ball of radius 0 a round a point $x \in \mathbb{R}^d$ by the symbol $B\varepsilon(x)$. We use the symbol \overline{B} to denote the closed unit

ball at the origin. For any set $X \subseteq \mathbb{R}^d$, the distance function and the projection map are defined by

$$\operatorname{dist}(x, X) := \inf_{y \in X} y - x$$
 and $P_X(x) := \underset{y \in X}{\operatorname{argmin}} y - x$,

respectively. Note that the function dist (\cdot, X) is 1-Lipschitz for any set $X \subseteq \mathbb{R}^d$, all $x \in \mathbb{R}^d$, and all $y \in P_X(x)$, we have

$$y - \overline{x} \le 2 x - \overline{x}.$$

We denote the diameter of a set X by

$$\operatorname{diam}(X) = \sup_{x, y \in X} x - y.$$

We call a function $h: \mathbb{R}^d \to \mathbb{R}$ *sublinear* if its epigraph is a closed convex cone, and in that case we define

$$Lin(h) := \{ x \in \mathbb{R}^d : h(x) = -h(-x) \}$$

to be its *lineality space*. Given a mapping $F: \mathbb{R}^d \to \mathbb{R}^m$ and a point $x \in \mathbb{R}^d$, we define

$$\lim_{\substack{x,x\to \overline{x}\\x=x}} \frac{F(x)-F(x)}{x-x}.$$

Given a mapping $F: \mathbb{R}^d \to \mathbb{R}^{m \times n}$ into the space of $m \times n$ matrices and a point $x \in \mathbb{R}^d$ then we define

$$\lim_{F} (\overline{x}) := \lim_{\substack{x,x \to \overline{x} \\ x = x}} \frac{F(x) - F(x)}{x - x},$$

where \cdot op denotes the operator norm defined on $\mathbb{R}^{m \times n}$.

Semialgebraicity. We call a set $X \subseteq \mathbb{R}^d$ *semialgebraic* if it is the union of finitely many sets defined by finitely many polynomial inequalities. Likewise, we call a function $f: \mathbb{R}^d \to \mathbb{R}$ semialgebraic if its graph gph $(f) = \{(x, f(x)) : x \in \mathbb{R}^d\}$ is semialgebraic.

Subdifferentials. Consider a locally Lipschitz function $f: \mathbb{R}^d \to \mathbb{R}^d$ and a point $x \in \mathbb{R}$. The Clarke subdifferential is the convex hull of limits of gradients evaluated at nearby points of differentiability:

$$\partial f(x) = \text{conv } \lim_{i \to \infty} \nabla f(x_i) : x_i \to x$$
,

where d is the set of points at which f is differentiable (recall Radamacher's theorem). If f is L-Lipschitz on a neighborhood U, then

for all
$$x \in U$$
 and $v \in \partial f(x)$, we have $v \le L$

This fact will be used throughout the paper. A point x satisfying $0 \in \partial f(x)$ is said to be critical for f. The Goldstein subdifferential, which appears in (1.2), will be a central object throughout. An important fact is that $\partial_{\sigma} f(x)$ is a closed convex set for any $x \in \mathbb{R}^d$ and $\sigma > 0$.

Manifolds. We will need a few basic results about smooth manifolds, which can be found in the references [10, 34]. A set $M \subseteq \mathbb{R}$ d is called a C p-smooth manifold around \overline{x} (with $p \ge 1$) if there exists a natural number m, an open neighborhood U of \overline{x} , and a C^p smooth mapping $F: U \to \mathbb{R}$ m such that the Jacobian $\nabla F(x)$ is surjective and $M \cap U = F^{-1}(0)$. The tangent and normal spaces to M at $x \in M$ near \overline{x} are defined to be $TM(x) = \ker(\nabla F(x))$ and $NM(x) = TM(x)^\perp = \operatorname{range}(\nabla F(x)^*)$, respectively. If M is a C^2 -smooth manifold around a point \overline{x} , then there exists $C \ge 0$ such that $y - x \in TM(x) + C$ y - x equiv 2B for all $x \cdot y \in M$ near \overline{x} . We also have that $x - PM(x) \in NM(PM(x))$ for all $x \in M$ near \overline{x} . Moreover, the projection mapping PM: $R^d \to R$ equiv B is C^{p-1} smooth on a neighborhood of \overline{x} and satisfies $\nabla PM(x) = P_{TM(x)}$ for all $x \in M$ near \overline{x} .

Covariant gradients and smooth extension Let $M \subseteq \mathbb{R}^d$ be a C^p -manifold around a point x for some $p \ge 1$. Then a function $f: M \to \mathbb{R}$ is called C^q -smooth (with $q \ge 1$) around the point x if there exists a C^q function $\hat{f}: U \to \mathbb{R}$ defined on an open neighborhood U of x and that agrees with f on $U \cap M$. In that case, the projection of $\nabla \hat{f}(x)$ onto $T_M(x)$ is independent of the choice of \hat{f} . We call this projection the covariant gradient of f at x and denote it by

$$\nabla_M f(x) := P_{T_M(x)}(\nabla \hat{f}(x)).$$

For example, the smooth extension

$$fM := f \circ PM$$

of f is $C^{\min\{p^{-1},q\}}$ smooth on a neighborhood of x and agrees with f along M. Thus, we will use the identification: $\nabla_M f(x) := \nabla f_M(x)$.

Active Manifolds. In this work, we will assume the local minimizer of interest lies on an *active manifold*. Informally, an active manifold is a smooth manifold along which the function varies smoothly and off of which the function varies sharply. We adopt the formal model of activity explicitly used in [26]. Related models exist, e.g., identifiable surfaces [52], manifolds of partial smoothness [38], VU-structures [36, 41], and $g \circ F$ decomposable functions [49].

Definition 1.2 (Active manifold) Consider a function $f : \mathbb{R}^d \to \mathbb{R}$ and fix a set $M \subseteq \mathbb{R}^d$ 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 neighborhood U of \overline{x} such that the following are true:

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

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

In the introduction, we provided two examples of functions that admit an active manifold at their minimizers. For example, function $f(u, v) = u^2 + |v|$ admits the active manifold $M := \{(u, 0) : u \in \mathbb{R}\}$ around the origin. On the other hand, the function f from (1.1) admits the active manifold $M = \{x \in \mathbb{R}^d : x_1 = x_2 = \dots x_m\}$ around the origin. To draw a distinction with partial smoothness property of [38], the function $f(x, y) = \max\{x, 0\} + y^2$ is partly smooth along the x axis, but the x-axis is not an active manifold for f around the origin; indeed, f does not satisfy the sharpness condition at the origin. We now turn to sublinear convergence guarantees.

2 Global Sublinear Convergent@exacent

The main goal of this work is to show that NTDescent locally converges nearly linearly for "typical" nonsmooth optimization problems. A natural question is whether NTDescent also possesses global nonasymptotic convergence guarantees. In this section, we prove two such guarantees: First, for arbitrary Lipschitz functions, we analyze the rate at which dist $(0, \partial_{\sigma_i} f(x_k))$ tends to zero as a function of k. Second, for convex Lipschitz functions, we analyze the rate at which $f(x_k)$ tends to inf f.

In the proofs of this section, the TDescent loop is ignored as we can only prove it terminates with descent near the minimizer. Instead, the global convergence guarantees follow from the properties on NDescent. Thus, our analysis essentially follows that of [21], where a nearly identical Minnorm method was introduced. The main difference between the NDescent and the method of [21] lies in the perturbation radius in Line 3 of Algorithm 1: while the radius of NDescent can be computed with access only to σ g_t , the radius in [21] requires knowledge of the Lipschitz constant of f, which we do not assume. Finally, we mention that [21] did not consider rates of convergence for convex problems.

Before stating the main result, we recall three key Lemmas, which underlie the proof. The first lemma shows that the vectors u_i and V_i generated by linesearch are Goldstein subgradients of decreasing norm.

Lemma 2.1 (Properties of linesearch) Let $f: \mathbb{R}^d \to \mathbb{R}$ be a locally Lipschitz function. Fix $x \in \mathbb{R}^d$, subgradient $g \in \partial f(x)$, budget T, and grid size G. Let u_i and V_i be generated by linesearch(x, g, G, T). Then

$$u_i, V_{i+1} \in \partial \sigma_i f(x)$$
 and $V_{i+1} \leq u_i \leq V_i$ (2.1)

for all $i = 0, \dots, G-1$.

Proof The proof follows by induction. We prove the base case only, since the induction is straightforward. First note that the inclusion $V_0 \in \partial f(x)$ implies that $u_0 \in \partial \sigma_0 f(x)$, since TDescent constructs u_0 as a convex combinations of subgradients evaluated in the ball $\overline{B}\sigma_0(\overline{x})$. Likewise, due to the argmin operation on line 4 of Algorithm 2, the subgradients generated by TDescent are decreasing in norm. Consequently, we have $u_0 \leq v_0$. A similar argument shows that $v_1 \in \partial \sigma_0 f(x)$ and $v_1 \leq u_0$. This completes the proof.

The next lemma shows that when f is convex, the minimal norm Goldstein subgradient may be used to bound the function values. We place the proof in Appendix A, since it follows from a standard argument.

Lemma 2.2 (Subgradient inequality) Suppose that $f: \mathbb{R}^d \to \mathbb{R}$ is a continuous convex function. Let $x, y \in \mathbb{R}^d$. Let L denote a Lipschitz constant for f on the ball $B_{2\sigma}(x)$. Then

$$f(x) - f(y) \le x - y \operatorname{dist}(0, \partial_{\sigma} f(x)) + 2\sigma L.$$

The final lemma provides conditions under which NDescent terminates with descent with high probability. The result is closely related to [21, Corollary 2.6], but we take extra care to analyze the perturbation radius in Line 3 of Algorithm 1.

Lemma 2.3 (NDescent loop terminates with descent) Let f be a locally Lipschitz function. Fix initial point $x \in \mathbb{R}^d$, radius $\sigma > 0$, subgradient $g \in \partial \sigma f(x)$, and failure probability $p \in (0, 1)$. Furthermore, let L be a Lipschitz constant of f on the ball $B_2\sigma(x)$. Suppose that

$$\sigma \leq \frac{\operatorname{dist}(0,\,\partial_{\sigma}\,f(x))}{\overline{128}L}; \quad and \quad T \geq \frac{64L^2}{\operatorname{dist}^2(0,\,\partial_{\sigma}\,f(x))} \quad 2\log(1/p) \ .$$

Define g_+ := NDescent (x, g, σ, T) . Then g_+ = 0 and the point x_+ := $x - \sigma - \frac{g_+}{g_+}$ satisfies

$$f(x_+) \le f(x) - \frac{\sigma_{\text{dist}}(0, \partial_{\sigma} f(x))}{8}$$
 with probability at least $1 - p$.

Proof First note that $g_+ \in \partial_\sigma f(x)$, so $g_+ \ge \operatorname{dist}(0, \partial_\sigma f(x)) > 0$. Now, observe that NDescent is precisely [21, Algorithm 1] with a different bound on the perturbation radius r. Indeed, in [21, Algorithm 1], r must satisfy

$$r < g_t$$
 $1 - 1 - \frac{g_t^2}{128L^2}$

for all $t \ge 0$. We now show that the constraint $r \le \sigma$ g_t implies the above bound. To that end, define the univariate function $h: a \to (1 - \frac{a^2}{128L^2})^2$. Then h

is increasing in a for $a \le L$. Moreover, for $a \in [0, L]$, we have $h(a) \ge \frac{\sqrt{a}}{128L}$. Consequently, since

$$\operatorname{dist}(0, \partial_{\sigma} f(x)) \leq g_t \leq L$$

for all $t \le T$, we have

$$r < \sigma \ g_t \leq \ \frac{\operatorname{dist}(0, \, \partial_\sigma \, f(x)) \ g_t}{\sqrt{128} L} \leq h(\operatorname{dist}(0, \, \partial_\sigma \, f(x))) \ g_t \leq \ h(\ g_t) \ g_t.$$

Thus the proof is a direct application of [21, Corollary 2.6].

Given these lemmata, we are now ready to state and prove our main sublinear convergence guarantee.

Theorem 2.4 (Sublinear convergence) Let $f: \mathbb{R}^d \to \mathbb{R}$ be a locally Lipschitz function. Fix initial point $x_0 \in \mathbb{R}^d$ and subgradient $g_0 \in \partial f(x_0)$. Assume that $g_0 = 0$. Let $L \in \mathbb{R} \cup \{+\infty\}$ be any Lipschitz constant of f over the widened sublevel set

$$S := \{ x + u : f(x) \le f(x_0) \text{ and } u \in \overline{B}(0) \}.$$

Fix a scalar $c_0 \in (0, 1]$, budget sequence $\{T_k\}$, grid size sequence $\{G_k\}$, and failure probability $p \in (0, 1)$. Let $\{x_k\}$ be generated by NTDescent $\{x, g, c_0, \{G_k\}, \{T_k\}\}$. Then for all $K \ge 0$, the following holds with probability at least 1 - p: Define $G := \min_{K \le k \le 2K - 1} G_k$ and $T := \min_{K \le k \le 2K - 1} T_k$. Then for all $i \le G$, the following bound holds with $\sigma_i := 2^{-(G-i)}$:

$$\min_{K \le k \le 2K - 1} \operatorname{dist}(0, \, \partial_{\sigma_i} f(x_k))$$

$$\le \max \quad \frac{8(f(x_K) - \inf f)}{\sigma_i K}, \, \frac{16L \quad \overline{2 \log(K G/p)}}{\overline{T}}, \, \sqrt[4]{\frac{128L}{T}} \sigma_i \quad .$$

Now suppose that f is convex and define $D := diam(\{x \in \mathbb{R}^d : f(x) \le f(x_0)\})$. Then

$$f(x_{2K-1}) - \inf f$$

$$\leq \min_{i \leq G} D \max \frac{8(f(x_K) - \inf f)}{\sigma_{i}K}, \frac{16L}{\sqrt[3]{T}}, \frac{2\log(KG/p)}{T}, \sqrt[4]{128L}\sigma_{i} + 2L\sigma_{i} .$$
(2.2)

Proof Let us assume that $L < +\infty$; otherwise the result is trivial. Fix K > 0 and $i \le G$. Define

$$_i := \max \quad \frac{16L \quad \overline{2 \log(K G' p)}}{\sqrt{\overline{T}}}, \quad \sqrt{128}L \sigma_i \quad .$$

For every $K \le k \le 2K - 1$, define

$$x_{k,i} := x_k - \sigma_i \frac{v_{i+1}}{v_{i+1}}, \quad \text{where } v_{i+1} := \text{NDescent}(x_k, u_i, \sigma_i, T_k),$$

and u_i appear in the definition of linesearch(x_k , g_k , max{ g_k , c_0 g_0 }, G_k , T_k); see Algorithm 3. Note that $V_{i+1} \in \partial \sigma_i f(x_k)$ by Lemma 2.1. Thus, in the event $\{\operatorname{dist}(0, \partial_{\sigma_i} f(x_k)) \geq i\}$, we have

- 1. $x_{k,i}$ is well-defined since $V_{i+1} = 0$;
- 2. the trust region constraint $\sigma_i \le \frac{v_{i+1}}{s}$ is satisfied for $s = \max\{g_k, c_0 g_0\}$ (in Algorithm 3); indeed,

$$\frac{v_{i+1}}{s} \geq \frac{\operatorname{dist}(0, \, \partial_{\sigma_i} f(x_k))}{s} \geq \frac{\sqrt{128L\sigma_i}}{s} \geq \sigma_i,$$

where the final inequality follows from the bound $s \le L$, a consequence of the inclusion $x_0 \subseteq \text{int } S$ and the Lipschitz continuity of f on S.

Finally, for every $K \le k \le 2K - 1$, define

$$A_{k,i} := f(x_{k,i}) - f(x_k) \ge -\frac{\sigma_i \operatorname{dist}(0, \, \partial_{\sigma_i} f(x_k))}{8} \cap \{\operatorname{dist}(0, \, \partial_{\sigma_i} f(x_k)) \ge i\}.$$

Now we apply Lemma 2.3.

To that end, observe that since $f(x_k)$ is nonincreasing and $\sigma_i \le 1/2$, every iterate x_k satisfies $B_2\sigma_i(x_k) \subseteq S$. Consequently, L is a Lipschitz constant of f on $B_2\sigma_i(x_k)$. Therefore, by Lemma 2.3, for every $K \le k \le 2K - 1$, we have

$$P(A_{k,i}) \le P(A_{k,i} \mid \operatorname{dist}(0, \partial_{\sigma_i} f(x_k)) \ge i) \le \frac{p}{GK}. \tag{2.3}$$

Thus, by a union bound, with probability at least $1 - \frac{p}{G}$, at least one of the following must hold at every index $K \le k \le 2K - 1$:

$$f(x_{k,i}) - f(x_k) \le -\frac{\sigma_i \operatorname{dist}(0, \, \partial_{\sigma_i} \, f(x_k))}{8} \quad \text{or} \quad \operatorname{dist}(0, \, \partial_{\sigma_i} \, f(x_k)) \le _i.$$

If $\operatorname{dist}(0, \partial_{\sigma_i} f(x_k)) \le i$ for some k satisfying $K \le k \le 2K - 1$, then the result follows. On the other hand, suppose that for all $K \le k \le 2K - 1$, we have $\operatorname{dist}(0, \partial_{\sigma_i} f(x_k)) > i$; in particular, we have $\operatorname{dist}(0, \partial_{\sigma_i} f(x_k)) > \overline{128}L\sigma_i$. Therefore, with probability at least $1 - \frac{p}{C}$, we must have

$$f(x_{k+1}) \leq f(x_k, i) \leq f(x_k) - \frac{\sigma_i \operatorname{dist}(0, \, \partial_{\sigma_i} \, f(x_k))}{8}, \quad \text{for all } K \leq k \leq 2K - 1.$$

where the first inequality follows since the trust region constraint is satisfied for $x_{k,i}$. Iterating this inequality, we have with probability at least $1 - \frac{p}{G}$, the bound

$$\min_{K \leq k \leq 2K-1} \operatorname{dist}(0, \, \partial_{\sigma_i} \, f(x_k)) \leq \, \frac{1}{K} \, \inf_{k=K} \operatorname{dist}(0, \, \partial_{\sigma_i} \, f(x_k)) \leq \, \frac{8(\, f(x_K) - \, f(x_{2K}))}{\sigma_i \, K}.$$

This proves the result for i. Taking a union bound over i then yields the bound for minimal norm Goldstein subgradient for all $i \le G$.

To prove (2.2), fix an $i \le G$ and let k_i be the index that attains the minimum. Then

$$f(x_{2K-1}) - \inf f \leq f(x_{k_i}) - \inf f \leq \operatorname{dist}(x_{k_i}, X_*) \min_{K \leq k \leq 2K-1} \operatorname{dist}(0, \partial_{\sigma_i} f(x_k)) + 2\sigma_i L$$

where the first inequality follows since $f(x_k)$ is nonincreasing and the second inequality follows from Lemma 2.2. The proof then follows from the upper bound $\operatorname{dist}(x_{k_i}, X) \leq D$.

theorem provides bounds on the minimal norm Goldstein subgradient within any window of indices $K \le k \le 2K - 1$. Let us briefly investigate the setting $T_k = k + 1$ for all $k \ge 0$. In this case, theorem implies that with probability at least 1 - p, we have

$$\begin{split} & \min_{K \leq k \leq 2K-1} \operatorname{dist}(0,\,\partial_{\sigma_i} \, f(x_k)) \\ & \leq \max \ \, \frac{8(\, f(x_K) - \inf \, f)}{\sigma_i \, K}, \, \frac{16L \quad \overline{2 \log(K \, G' \, p)}}{2\overline{K}}, \, \sqrt[4]{128}L \, \sigma_i \end{split}$$

for all $i \le G$. Let us now suppose G is large enough that there exists $i \le G$ satisfying $(1/2)K^{-1/2} \le \sigma_i \le K^{-1/2}$, e.g., we may assume $G_k = \log(k^{1/2})$ for all $k \ge 0$. Then, we find that at most $O(KTG) = O(K^2G)$ first-order oracle evaluations are needed to find a point x_k satisfying

$$dist(0, \partial_{K^{-1/2}} f(x_k)) = \tilde{O}(K^{-1/2}),$$

where \tilde{O} hides logarithmic terms in G, K and p. Let's consider two settings for G_k.

- 1. **Setting 1**: $G_k = O(\log(k^{1/2}))$. In this case, NTDescent finds a point x_k satisfying
 - $\operatorname{dist}(0,\ \partial_{\varepsilon}\,f(x_k)) \leq \varepsilon \ \text{ using at most } \tilde{O}(\varepsilon^{-4}) \text{ first-order oracle evaluations}.$
- 2. **Setting 2**: $G_k = k + 1$. In this case, NTDescent finds a point x_k satisfying $\operatorname{dist}(0, \partial_{\varepsilon} f(x_k)) \leq \varepsilon$ using at most $\tilde{O}(\varepsilon^{-6})$ first-order oracle evaluations.

The complexity of Setting 1 is more favorable than the complexity of Setting 2. Nevertheless, when we establish our local rapid convergence guarantees, we will work in Setting 2, which has more favorable local convergence properties. Before moving on, we note that the above guarantees likewise apply in the convex setting, namely

NTDescent finds a point x_k with $f(x_k) - f^* \le \varepsilon$ using at most $\tilde{O}(\varepsilon^{-4})$, respectively $\tilde{O}(\varepsilon^{-6})$, first-order oracle evaluations in Setting 1, respectively Setting 2.

In addition to the nonasymptotic guarantees of Theorem 2.4, the reader may wonder whether a given limit point \overline{x} of NTDescent is Clarke critical, meaning $0 \in \partial f(\overline{x})$. We prove that this is indeed the case under a bounded sublevel set condition. We place the proof in Appendix C since it follows a similar line of reasoning as Theorem 2.4.

Corollary 2.5 (Limiting points are Clarke critical) Let $f: \mathbb{R}^d \to \mathbb{R}$ be a locally Lipschitz function. Fix initial point $x_0 \in \mathbb{R}^d$ and subgradient $g_0 \in \partial f(x_0)$. Assume that $g_0 = 0$. Suppose the sublevel set $\{x: f(x) \le f(x_0)\}$ is bounded. Fix scalar $c_0 \in (0, 1]$, budget sequence $\{T_k\}$, grid size sequence $\{G_k\}$ such that $\{G_k\}$ tends to infinity and $T_k \ge k$. Let $\{x_k\}$ be generated by NTDescent $\{x, g, c_0, \{G_k\}, \{T_k\}\}$. Then with probability one, all the limiting points of $\{x_k\}$ are Clarke critical.

This concludes our sublinear convergence guarantees for NTDescent. In the following section, we describe the key structural assumptions needed to ensure that NTDescent locally rapidly converges.

3 Main Assumption, Examples, and Consequences

In this section, we introduce our key structural assumption—Assumption A. In Sect. 3.1 we show that Assumption A holds for generic semialgebraic functions and certain properly C^p decomposable functions. Then, in Sect. 3.2, we extract several key consequences of Assumption A. These consequences will be instrumental in proving the gradient inequality (1.5) and rapid convergence of NTDescent. We now turn to the assumption.

Assumption A Function $f: \mathbb{R}^d \to \mathbb{R}$ is locally Lipschitz with local minimizer $\bar{x} \in \mathbb{R}^d$.

(A1) **(Quadratic Growth)** There exists Y > 0 such that

$$f(x) - f(\bar{x}) \ge \frac{y}{2} |x - \bar{x}|^2$$
 for all x near \bar{x} .

- (A2) (Active Manifold) Function f admits a C^4 -smooth active manifold M around \overline{x} .
- (A3) **(Strong-**(*a*) **regularity)** There exists C(a) > 0 such that

$$P_{T_M (y)}(v - \nabla_M f(y)) \le C_{(a)} x - y$$

for all $x \in \mathbb{R}^d$, $v \in \partial f(x)$, and $y \in M$ near \bar{x} .

(A4) $((b \le)$ -regularity) The following inequality holds

$$f(y) \ge f(x) + V, y - x + o(y - x)$$

as $y \to x$ and $x \to x$ with $v \in \partial f(x)$,

where $o(\cdot)$ is any univariate function satisfying $\lim_{t\to 0} o(t)/t = 0$.

Some comments are in order. Assumption (A1) is a classical regularity condition that ensures local linear convergence of gradient methods for smooth convex functions. Assumptions (A2), (A3), and (A4) describe the interaction of f and a distinguished smooth manifold M. Assumption (A2) requires M to be an active manifold for faround \bar{x} in the sense of Definition 1.2. In particular, along the manifold M, the function f is C^4 smooth with covariant gradient $\nabla_M f$; see Sect. 1.7 for a definition. Assumption (A3) shows that in tangent directions the covariant gradient along the manifold approximates the subgradients of f up to a linear error. This property recently appeared in [5, 19], where it was used to study saddle avoidance properties of the subgradient method for nonsmooth optimization. Finally, Assumption (A4) is a restricted lower smoothness property, showing that linear models of foff the manifold are underapproximators of f on the manifold up to first-order. Note that the property is automatic if f is weakly convex, meaning the mapping $x \to f(x) + \frac{\rho}{2} x^2$ is convex for some $\rho \ge 0$. The weakly convex class is broad and contains all compositions of convex functions with smooth mappings that have Lipschitz Jacobians; see the survey [17] for an introduction. We mention that the name $(b \le)$ -regularity" is motivated by "uniform semismoothness" property of [19], which was called the "(b)-regularity property."

In the following section, we provide examples of functions satisfying Assumption A.

3.1 Examples of Assumption A

In this section, we show that the aforementioned problems satisfy Assumption A. The most important example is the class of generic semialgebraic functions. The following theorem is essentially contained in [19, 23], but we provide a proof for completeness.

Theorem 3.1 (Generic semialgebraic functions) Consider a locally Lipschitz semial-gebraic function $f: \mathbb{R}^d \to \mathbb{R}$. Then for a full Lebesgue measure set of $w \in \mathbb{R}^d$, the tilted function $fw: x \to f(x)+w$ x satisfies Assumption A at every local minimizer.

Proof The proof is a consequence of [19, Theorem 3.31] and [23, Corollary 4.8, Theorem 4.16]. A combination of Corollary 4.8 and Theorem 4.16 in [23] shows that for a full Lebesgue measure set of $W \in \mathbb{R}^d$, the following hold: every local minimizer \overline{x} of f_W lies on a C^4 active manifold M, verifying (A2); and the quadratic growth condition (A1) holds at \overline{x} . Next, [19, Theorem 3.31] shows that f_W also satisfies the strong (a) property (A3) along M; applying [19, Theorem 3.11 and Theorem 3.4], we deduce that f_W also satisfies the (b_S)-regularity property (A4) along M at \overline{x} .

Turning to our second class, we introduce so-called *properly C* ^p *decomposable* functions, originally proposed and analyzed in [49]. At a high-level, the class consists of functions that are locally the composition of a sublinear function with a smooth mapping, which together satisfy a transversality condition.

Definition 3.2 (Decomposable functions) A function $f: \mathbb{R}^d \to \mathbb{R}$ is called *properly* C^p *decomposable at* \overline{x} *as* $h \circ c$ if near \overline{x} it can be written as

$$f(x) = f(\bar{x}) + h(c(x))$$

for some C^p -smooth mapping $c: \mathbb{R}^d \to \mathbb{R}^m$ satisfying $c(\overline{x}) = 0$ and some proper, closed sublinear function $h: \mathbb{R}^m \to \mathbb{R}$ satisfying the transversality condition:

$$\lim(h) + \operatorname{range}(\nabla c(\bar{x})) = R^m.$$

The following theorem shows that decomposable functions satisfy Assumption A near local minimizers if they also satisfy a strict complementarity condition and a quadratic growth bound. The proof is a consequence of results found in works [19, 26, 38, 49].

Theorem 3.3 (Properly decomposable functions) Consider a locally Lipschitz function $f: \mathbb{R}^d \to \mathbb{R}$. Let \bar{x} be a local minimizer of f and suppose that f is properly C^4 decomposable at \bar{x} . Furthermore, suppose that

- 1. (Strict Complementarity) We have that $0 \in \operatorname{ri} \partial f(\overline{x})$.
- 2. (Quadratic growth) There exists Y > 0 such that

$$f(x) - f(\bar{x}) \ge \frac{\gamma}{2} |x - \bar{x}|^2$$
 for all x near \bar{x} .

Then f satisfies Assumption A at \bar{x} .

Proof To set the notation for the proof, recall that since f is properly C decomposable, there exist functions h and c satisfying the conditions of Definition 3.2. The discussion in [49, p. 683-4] then shows that the set

$$M := c^{-1}(\ln(h))$$

is a so-called C^4 manifold of partial smoothness for f around \overline{x} in the sense of Lewis [38]. Moreover, f is prox-regular at \overline{x} for 0 in the sense of [46, Definition 1.1], since by definition it is strongly amenable [46, Definition 2.4] at; see [46, Proposition 2.5]. Thus, according to [31, Theorem 5.3], partial smoothness, prox-regularity, and strict complementarity ensure that the sharpness condition of Definition 1.2 holds. Consequently, M is a C^4 smooth active manifold around \overline{x} , verifying (A2). In addition, [19, Corollary 3.24] ensures that f satisfies the (A3) and (A4) properties along M.

A popular class of decomposable objectives arises from pointwise maxima of smooth functions that satisfy an affine independence property. For example, this class was considered in the work of Han and Lewis [30]. As an immediate corollary of Theorem 3.3, we show that such functions satisfy Assumption A.

Corollary 3.4 (Max-of-smooth functions) Consider a locally Lipschitz function f and a family of C^4 smooth functions $f_i : \mathbb{R}^d \to \mathbb{R}$ indexed by a finite set $i \in I$. Fix a local

minimizer \bar{x} of f and suppose the set $\{\nabla f_i(\bar{x})\}_{i\in I}$ is affinely independent. Suppose furthermore that f is locally expressible as

$$f(x) := \max_{i \in I} f_i(x)$$
 for all x near x .

Then provided the strict complementarity and quadratic growth conditions of Theorem 3.3 hold, the function f satisfies Assumption A at \overline{x} .

Proof To prove the result, note that the affine independence property is simply a restatement of the transversality condition of Definition 3.2 for the smooth mapping $x \to (f_i(x))_{i \in I}$ and the sublinear function $y \to \max_{i \in I} y_i$.

We now turn our attention to the key consequences of Assumption A.

3.2 Key Consequences of Assumption A

The following proposition summarizes the key consequences of Assumption A. The proof of the result is straightforward but technical, so we place it in Appendix B.

Proposition 3.5 (Consequences of Assumption A) Suppose f satisfies Assumption A at \overline{x} . Then there exists $\delta_A > 0$ such that on the ball $B_2\delta_A(\overline{x})$, the projection operator PM is C^3 with Lipschitz Jacobian and the smooth extension $fM := f \circ PM$ is C^3 with Lipschitz gradient. Moreover, the following bounds hold:

- 1. **(Quadratic growth)** The quadratic growth bound (A1) holds throughout $B_{\delta_A}(\vec{x})$.
- 2. (Smoothness of PM) For all $x \in B_{\delta_A}(x)$ and $x \in B_2\delta_A(x)$, we have

$$P_M(x) - P_M(x) - P_{T_M(P_M(x))}(x - x) \le C_M(\operatorname{dist}^2(x, M) + x - x^{-2}),$$
(3.1)

where $C_M := 2 \operatorname{lip}_{\nabla P_M}^{\operatorname{op}} (\vec{x})$.

3. **(Bounds on** $\nabla_M f$) For all $x \in B_{\delta_A}(x)$, we have

$$\frac{y}{2} P_M(x) - \overline{x} \le \nabla M f(P_M(x)) \le \beta P_M(x) - \overline{x}, \qquad (3.2)$$

where $\beta := 2 \operatorname{lip}_{\nabla f_M}(\vec{x})$.

4. (Consequence of strong (a)) For all $x \in B_{\delta_A}(x)$ and $\sigma \le \delta_A$, we have

$$\sup_{q \in \partial \sigma} P_{TM} \left(P_{M} \left(x \right) \right) \left(g - \nabla_{M} f \left(P_{M} \left(x \right) \right) \right) \le C_{(a)} \left(\operatorname{dist}(x, M) + \sigma \right); \tag{3.3}$$

$$\sup_{g \in \partial \sigma \, f(x)} P_{TM \, (PM \, (x))} g \leq C_{(a)} (\operatorname{dist}(x, M) + \sigma)$$

$$+ \beta P_M(x) - \bar{x};$$
 (3.4)

$$\sup_{g \cdot g \in \partial \sigma f(x)} P_{T_M(P_M(x))}(g - g) \le 2C_{(a)}(\operatorname{dist}(x, M) + \sigma). \quad (3.5)$$

5. (Aiming) For all $x \in B\delta_A(\vec{x})$ and all $v \in \partial f(x)$, we have

$$V, x - P_M(x) \ge \mu \operatorname{dist}(x, M), \tag{3.6}$$

where $\mu := \frac{1}{4} \liminf_{\substack{x \to -\infty \\ x \to -x}} \mathrm{dist}(0, \partial_f(x)).$ 6. (Subgradient bound) For all $x \in B\delta_{\mathrm{A}}(\vec{x})$ and $\sigma \leq \delta_{\mathrm{A}}$, we have

$$\sup_{g\in\partial\sigma\;f(x)}\;g\leq\;L,$$

where $L := 2 \operatorname{lip}_{f}(\vec{x})$ 7. **(Function gap)** For all $x \in B_{\delta_A}(x)$, we have

$$f(x) - f(\bar{x}) \le L \operatorname{dist}(x, M) + \frac{\beta}{2} P_M(x) - \bar{x}^2.$$
 (3.7)

Let us briefly comment on the result. Item 2 provides a crucial smoothness property of the projection operator of M. Item 3 shows that the Riemannian gradient of f is proportional to the distance of the projection y to \bar{x} . Item 4 shows how the Goldstein subgradients inherit the strong (a) property (A3) of Assumption A. Indeed, Equation (3.4) shows that Goldstein subgradients are "small" in tangent directions and Equation (3.5) shows Goldstein subgradients vary in an approximate Lipschitz fashion in tangent directions. Item 5 shows that the subgradients of *f* off of the manifold have a constant level of correlation with x - PM(x), i.e., the direction -v "aims" towards the manifold. Note that $\mu > 0$ due to the active manifold Assumption (A2). The proof of Item 5 is based on Assumptions (A2) and (A4); a similar result appears in [18, Theorem D.2]. Item 6 provides a bound on the Goldstein subgradients of we will appeal to this bound many times throughout the analysis without referencing this proposition. Finally, Item 7 decomposes the function gap into a sum of two terms: the distance to the manifold and the squared distance of the projection to the solution. The proof relies on the smoothness of *f* along the manifold. Note that the trivial upper bound $L \times \overline{x}$ for the gap can be weaker than (3.7).

This concludes our discussion of Assumption A. The following three sections establish further consequences: the gradient inequality (1.5) (Sect. 4); rapid local convergence of NDescent and TDescent (Sect. 5); and rapid local convergence of NTDescent (Sect. 6). In all three sections, we use the notation and results introduced in Proposition 3.5. Finally, the statements of the results in Sect. 4 and 5 contain several parameters/radii which we will use in Sect. 6 to determine the region of near linear convergence and the oracle complexity for NTDescent. For the readers' convenience, we have listed these parameters in Table 1.

Table 1	Parameters used
throughout Sects. 4 and 5	

Parameter	Definition
$\overline{D_1}$	μ 8(μ+ 1.) μ 2 γ 4
D_2	$\frac{\mu}{2}$
C_1	$\frac{y}{4}$
C_2	$\min \frac{\gamma}{8C(a)}, \ \frac{\min\{1, 1/\delta_{\text{A}}\}}{2}$
C_3	$\frac{C_1^2}{8L}$
C_4	min $\frac{\beta}{C_{(a)}(1+\delta_{A})}$, $\frac{\min\{\mu/\delta_{A}, C_{3}D_{2}/\beta\}}{4(1+(1+\delta_{A})C_{M})(\mu+L))}$, $\frac{1}{2}$
C_5	min $\frac{\beta}{2C(a)}$, $\frac{C_3D_2}{32C(a)\beta}$, C_4 , $\frac{C_2}{4}$
$\delta_{ m GI}$	min $\frac{\delta_A}{4}$, $\frac{D_1}{CM}$
δ_{ND}	min δ_{GI} , $\frac{D_2}{D_1L}$ $\frac{D_2}{128}$
$\delta_{ m Grid}$	min $\frac{\delta_{A}}{2}$, $\frac{1}{CM (D_{1}^{-1}+1)}$, $\frac{\mu}{8(C(a)+\beta)}$

4 Verifying the Gradient Inequality (1.5) Under Assumption A

In this section, we establish the gradient inequality (1.5) for functions satisfying Assumption A. Throughout the section, we assume that Assumption A is in force. We also use the notation set out in Proposition 3.5.

We present the formal statement and the gradient inequality (1.5) in Theorem 4.3, which appears at the end of this section. The proof is a consequence of the two lemmata. In the first lemma, we prove a constant-sized lower bound for $\operatorname{dist}(0, \partial_{\sigma} f(x))$, whenever σ is sufficiently small. The proof of this bound relies on the active manifold assumption (A2) and the aiming inequality (3.6). A consequence of the argument is that all elements of $\partial_{\sigma} f(x)$ are correlated with the normal direction $x - PM(x) \in NM(PM(x))$. Later in Proposition 5.1 we will also show that Algorithm 1 (NDescent) terminates rapidly when σ is in the region, motivating the name Normal Descent. We now turn to the lemma.

Lemma 4.1 (Lower bound on Goldstein subgradients I) Define

$$D_1:=\frac{\mu}{8(\mu+L)}; \qquad D_2:=\frac{\mu}{2}; \qquad and \qquad \delta_{\rm GI}:=\min \ \frac{\delta_{\rm A}}{4}, \frac{D_1}{CM} \ .$$

Then for all $x \in B\delta_{G}(x)$ and $0 < \sigma \le D_1 \mathrm{dist}(x, M)$, we have

123

$$\mathrm{dist}(0,\,\partial_\sigma\,f(x))\geq\ D_2.$$

Proof We begin with some preliminary bounds. Fix $x \in B\delta_{GI}(\vec{x})$ and $\sigma > 0$ satisfying the lemma assumptions. We observe that

$$\sigma \leq D_1 \operatorname{dist}(x, M) \leq \operatorname{dist}(x, M) \leq x - \overline{x} \leq \delta_{GI},$$

where the second inequality follows since $D_1 \le 1$ and the third follows since $x \in M$. Consequently,

$$\begin{split} LC_M\left(\sigma^2 + \operatorname{dist}^2(x,\,M)\right) &\leq \delta_{\mathrm{GI}} LC_M\left(\sigma + \operatorname{dist}(x,\,M)\right) \\ &\leq 2L\,\delta_{\mathrm{GI}} C_M\,\operatorname{dist}(x,\,M) \\ &\leq 2L\,D_1 \mathrm{dist}(x,\,M), \end{split} \tag{4.1}$$

where the first inequality follows from the bound $\max\{\sigma, \operatorname{dist}(x, M)\} \le \delta$ GI and the second follows from the bound $\sigma \le \operatorname{dist}(x, M)$. We now turn to the proof.

Now, let $x \in \overline{B}\sigma(x) \subseteq B\delta_A(x)$ and observe that by aiming condition (3.6),

$$V, x - P_M(x) \ge \mu \operatorname{dist}(x, M)$$
 for all $v \in \partial f(x)$.

We claim that $V, x - PM(x) \ge D_2 \operatorname{dist}(x, M)$ for all $v \in \partial f(x)$. Indeed, for all $v \in \partial f(x)$ we may upper bound the inner product as follows:

$$V, x - P_M(x)$$

 $\leq V, x - P_M(x) + V \quad x - P_M(x) - x + P_M(x)$
 $\leq V, x - P_M(x) + L(I - P_{T_M(P_M(x))})(x - x) + LC_M(\sigma^2 + \text{dist}^2(x, M))$
 $\leq V, x - P_M(x) + 3L D_1 \text{dist}(x, M),$

where the second inequality follows from the bound $v \le L$ and Item 2 of Proposition 3.5; and the third inequality follows from $x - x \le \sigma \le D_1 \text{dist}(x, M)$ and (4.1). Consequently, for all $v \in \partial f(x)$, we have

$$V, x - P_M(x) \ge \mu \operatorname{dist}(x, M) - 3L D_1 \operatorname{dist}(x, M)$$

$$\ge \mu \operatorname{dist}(x, M) - \mu \sigma - 3L D_1 \operatorname{dist}(x, M)$$

$$\ge \mu (1 - D_1(1 + 3L/\mu)) \operatorname{dist}(x, M)$$

$$= D_2 \operatorname{dist}(x, M), \tag{4.2}$$

where the second inequality follows from 1-Lipschitz continuity of dist (\cdot, M) ; and the final inequality follows from the bound $D_1 \le \frac{1}{2(1+3L/U)}$. This proves the claim.

Now, fix $g \in \partial \sigma f(x)$. By definition of $\partial \sigma f(x)$, there exists a family of coefficients $\lambda_i \in [0, 1]$, points $x_i \in \overline{B}\sigma(x) \subseteq B\delta_A(x)$, and subgradients $g_i \in \partial f(x_i)$ indexed by a finite set $i \in I$ such that $i \in I$ and $g = \int_{i \in I} \lambda_i g_i$. Thus, by (4.2), we have

$$g, x - P_M(x) = \lambda_i g_i, x - P_M(x) \ge D_2 \operatorname{dist}(x, M).$$

Therefore, we have

$$g \ge \frac{g, x - P_M(x)}{\operatorname{dist}(x, M)} \ge D_2,$$

as desired.

In the second lemma, we provide a lower bound for dist $(0, \partial_{\sigma} f(x))$ on the order of $PM(x) - \overline{x}$, provided $\sigma = O(PM(x) - \overline{x})$. The proof of this bound relies on quadratic growth (A1) and strong (a)-regularity (A3). A consequence of the argument is that the minimal norm element of $\partial_{\sigma} f(x)$ is close to the tangent vector $\nabla_M f(PM(x)) \in TM(PM(x))$. Later in Proposition 5.6 we will also show that Algorithm 2 (TDescent) terminates rapidly when σ is in the region, motivating the name Tangent Descent. We now turn to the lemma.

Lemma 4.2 (Lower bound on Goldstein subgradients II) Define

$$C_1:=rac{\gamma}{4}; \quad and \quad C_2:=\min \ rac{\gamma}{8C_{(a)}}, rac{\min\{1,\ 1/\delta_{\mathrm{A}}\}}{2} \ .$$

Then for all $x \in B_{\delta_A}(x)$ and $\sigma \ge 0$ satisfying

$$\max\{\operatorname{dist}(x, M), \sigma\} \leq C_2 P_M(x) - \overline{x},$$

we have

$$P_{T_M(P_M(x))}(g) \ge C_1 P_M(x) - \overline{x}$$
 for all $g \in \partial \sigma f(x)$.

Proof For the purposes of this proof, the term $1/\delta_A$ in the definition of C_2 is unnecessary; however, it will be crucial in the proof of Theorem 4.3. Turning to the proof, fix $x \in B\delta_A(x)$ and $\sigma \ge 0$ satisfying the lemma assumptions. Define y = PM(x). Note that

$$\sigma \le C_2 \ y - \overline{x} \le 2C_2 \ x - \overline{x} \le \delta_A$$

Thus, by (3.3), for all $g \in \partial \sigma f(x)$, we have

$$P_{T_M(y)}(g - \nabla_M f(y)) \le C_{(a)}(\operatorname{dist}(x, M) + \sigma) \le \frac{\gamma}{4} y - \overline{x}.$$

In addition, by (3.2), we have $\nabla M f(y) \ge \frac{y}{2} y - \overline{x}$. Therefore, for all $g \in \partial_{\sigma} f(x)$, we have

$$P_{T_M(y)}(g) \geq \nabla \quad M \ f(y) - \ C_{(a)}(\operatorname{dist}(x,M) + \sigma) \geq \ \frac{\gamma}{4} \ y - \bar{x},$$

as desired.

Given these lemmata, we are now ready to establish the gradient inequality (1.5). The following theorem verifies the bound

$$\sigma_{\mathrm{dist}}(0, \, \partial_{\sigma} f(x)) \geq \eta(f(x) - f(x)),$$

for some $\eta > 0$ provided x is sufficiently near \bar{x} and (x, σ) lies within one of two regions, described in Item 1 and Item 2 of Theorem 4.3. Item 1 and Item 2 roughly correspond to the regions considered in Lemma 4.1 and Lemma 4.2, Comparing with the statement of the gradient inequality (1.5), we see that gradient inequality of Theorem 4.3 does not require knowledge of an explicit function $\sigma(x)$. Instead, we need only find some σ proportional to D_1 dist(x, M) or C_2 P_M (x) – \overline{x} up to a factor of, say, 2. Later in Proposition 6.1 we show that this flexibility allows us to find an appropriate σ through the linesearch procedure.

Theorem 4.3 (Gradient inequality) Suppose that function f satisfies Assumption A at $\overline{x} \in \mathbb{R}^d$. For any constants $a_1 \in (0, D_1]$ and $a_2 \in (0, C_2]$, we have

$$\sigma_{\mathrm{dist}}(0,\,\partial_{\sigma}\,f(x))\geq \min \ \frac{y\,a_2}{8\,\max\{4La_2^2,\,\beta\}}, \frac{\mu_{a_1}}{4\,\max\{2L\,,\,\beta/a_2^2\}} \ (f(x)-f(x)),$$

whenever $x \in B\delta_{C1}(x)$ and $\sigma > 0$ satisfy Item 1 or Item 2:

- 1. (a) $\frac{a_1}{2} \operatorname{dist}(x, M) \le \sigma \le a_1 \operatorname{dist}(x, M)$; (b) $a_2^2 \ PM(x) \frac{x}{x}^2 \le \operatorname{dist}(x, M)$. 2. (a) $\frac{a_2}{2} \ PM(x) x \le \sigma \le a_2 \ PM(x) x$;
 - (b) $\frac{\text{dist}(x, M)}{\sigma} \le 2a_2 \ PM(x) x$.

Moreover, for any $x \in B\delta_{CI}(\overline{x})\setminus \{\overline{x}\}$, there exists $\sigma > 0$ such that Item 1 or Item 2 is satisfied.

Proof We first show that for any $x \in B\delta_{GL}(\overline{x})\setminus\{\overline{x}\}$, there exists $\sigma > 0$ such that either Item 1 or Item 2 is satisfied. We consider two cases. First, suppose $x \in M$. Then Item 2 is trivially satisfied for $\sigma = a_2 PM(x) - x$. Second, suppose $x \notin M$ and Item 1 cannot be satisfied for any $\sigma > 0$. In this case, we have

$$dist(x, M) \le a_2^2 P_M(x) - \bar{x}^2 = 2a_2\sigma P_M(x) - \bar{x}$$
 with $\sigma := a_2 P_M(x) - \bar{x}/2$.

Thus, Item 2 is satisfied.

Now we prove the gradient inequality is satisfied whenever σ satisfies Item 1 or Item 2. Let us suppose that Item 1 holds for some $x \in B\delta_{G}(x)$ and $\sigma > 0$. From (3.7), we have the bound:

$$\begin{split} \frac{1}{\max\{2L,\,\beta/a_2^2\}} (f(x) - f(\bar{x})) &\leq \frac{1}{\max\{2L,\,\beta/a_2^2\}} \quad L \mathrm{dist}(x,\,M) + \frac{\beta}{2} \;\; P_M(x) - \bar{x}^{\;\;2} \\ &\leq \frac{1}{2} \;\; \mathrm{dist}(x,\,M) + \;\; a_2^2 \;\; P_M(x) - \bar{x}^{\;\;2} \\ &\leq \mathrm{dist}(x,\,M). \end{split}$$

Now observe that the assumptions of Lemma 4.1 are satisfied since $x \in B\delta_{GI}(\vec{x})$, $a_1 \le D_1$, and x and σ satisfy Item 1. Therefore, we have

$$\sigma_{\mathrm{dist}}(0,\,\partial_{\sigma}\,f(x))\geq\sigma\,D_{2}\geq\frac{\mu\,a_{1}}{4}\mathrm{dist}(x,\,M)\geq\,\frac{\mu\,a_{1}}{4\,\max\{2L\,,\,\beta/\,a_{2}^{2}\}}(f(x)-\,f(x)),$$

as desired.

Next, let us suppose that Item 2 holds for some $x \in B\delta_{GI}(\vec{x})$ and $\sigma > 0$. From (3.7), we have the bound:

$$\frac{1}{\max\{4La_{2}^{2},\beta\}}(f(x)-f(\overline{x})) \leq \frac{1}{\max\{4La_{2}^{2},\beta\}} L \operatorname{dist}(x,M) + \frac{\beta}{2} P_{M}(x) - \overline{x}^{2}$$

$$\leq \frac{1}{2} \frac{\operatorname{dist}(x,M)}{2a_{2}^{2}} + P_{M}(x) - \overline{x}^{2}$$

$$\leq \frac{1}{2} \frac{\operatorname{dist}(x,M) P_{M}(x) - \overline{x}^{2}}{2a_{2}\sigma} + P_{M}(x) - \overline{x}^{2}$$

$$\leq P_{M}(x) - \overline{x}^{2}.$$

Now observe that since $a_2 \le C_2$ and x and σ satisfy Item 2, we have

$$\sigma \le C_2 P_M(x) - \overline{x} \le 2C_2 \delta_{GI} \le (1/\delta_A)(\delta_A/4) \le 1$$

where we use the bound $C_2 \le 1/2\delta_A$. Consequently, we have

$$\operatorname{dist}(x, M) \le 2\sigma C_2 \ P_M(x) - \overline{x} \le C_2 \ P_M(x) - \overline{x}.$$

Therefore, $\max\{\operatorname{dist}(x, M), \sigma\} \le C_2 \ PM(x) - \overline{x}$, so the conditions of Lemma 4.2 are satisfied (recall $\delta_{\mathrm{GI}} \le \delta_{\mathrm{A}}$). Thus, let g denote the minimal norm element of $\mathfrak{P}_{\sigma} f(x)$ and let us apply Lemma 4.2:

$$\operatorname{dist}(0,\,\partial_{\sigma}\,f(x)) = \quad g \geq \quad P_{TM\ (PM\ (x))}(g) \geq \quad \frac{\gamma}{4} \ P_{M}\left(x\right) - \bar{x}.$$

Consequently, we have

$$\sigma_{\mathrm{dist}}(0,\,\partial_{\sigma}\,f(x))\geq\,\frac{\sigma\,\gamma}{4}\;\;P_{M}\left(x\right)-\frac{1}{x}\geq\,\,\frac{\gamma\,a_{2}}{8\,\max\{4La_{2}^{\,2},\,\beta\}}(\,f(x)-\,f(x)),$$

where the last inequality follows from $\sigma \geq \frac{a_2}{2} P_M(x) - \overline{x}$. This completes the proof.

Remark 1 Note that a_1 , $a_2 \in (0, 1)$ as claimed in Sect. 1.3.1, where we introduced the *normal and tangent regions* appearing in the statement of Theorem 4.3.

This concludes the proof of the gradient inequality (1.5) under Assumption A. In Sect. 6, we will use the gradient inequality to establish rapid local convergence of NTDescent. Before proving that, the following section analyzes TDescent and NDescent methods.

5 Rapid Termination Defision and TDescent Under Assumption A

In this section, we analyze the NDescent and TDescent methods, showing that both methods rapidly terminate with descent in appropriate regions. Throughout the section, we assume that Assumption A is in force. We also use the results and notation of Proposition 3.5, Table 1, Lemma 4.1, and Lemma 4.2.

The main results of this section are Propositions 4.1 and 5.6, which analyze NDescent and TDescent, respectively. Proposition 5.1 shows tha NDescent terminates with descent in a constant number of iterations within the region considered in Item 1 of Theorem 4.3. Proposition 5.6 shows that TDescent either terminates with descent in $O(\log^{-1}(f(x) - f(x)))$ iterations or f(x) - f(x) is already exponentially small in T within the region considered in Item 2 of Theorem 4.3. These lemmata will be the basis of our main convergence theorem—Theorem 6.3—appearing in Sect. 6.

5.1 Analysis of NDescent

The following proposition shows that NDescent locally terminates in finitely many iterations whenever σ is sufficiently small. The result is a simple consequence of Lemmas 2.3 and 4.1.

Proposition 5.1 (NDescent loop terminates with descent) Define a radius

$$\delta_{\text{ND}} := \min \ \delta_{\text{GI}}, \frac{D_2}{D_1 L} \ \overline{128} \ .$$

Then for all $x \in B\delta_{ND}(\vec{x})$, radii $\sigma > 0$ with $\sigma \leq D_1 \text{dist}(x, M)$, subgradients $g \in \partial_{\sigma} f(x)$, failure probabilities $p \in (0, 1)$ and budgets T > 0 satisfying

$$T \ge \frac{64L^2}{D_2^2} - 2\log(1/p)$$
,

the point $x_+ := NDescent(x, g, \sigma, T)$ satisfies

$$f(x_+) \le f(x) - \frac{\sigma_{\text{dist}}(0, \, \partial_{\sigma} f(x))}{8}$$
 with probability at least $1 - p$.

Proof Fix $x \in B_{\delta_{ND}}(\overline{x})$ and $\sigma > 0$ satisfying the lemma assumptions. Observe that

$$\sigma \le D_1 \operatorname{dist}(x, M) \le D_1 \delta_{ND} \le \min \delta_{GI}, \frac{D_2}{L 128}$$

where the final inequality follows from the bound $D_1 \le 1$; see Lemma E.1. Thus, by Lemma 4.1, we have dist $(0, \partial_{\sigma} f(x)) \ge D_2$ (recall $\delta_{ND} \le \delta_{GI}$). Consequently,

$$\sigma \leq \frac{D_2}{L} \leq \frac{\operatorname{dist}(0, \partial_{\sigma} f(x))}{L}.$$

Therefore, σ and T satisfy the assumptions of Lemma 2.3. Hence, the desired descent condition is guaranteed with probability at least 1 - p.

We now turn to the analysis of the TDescent step.

5.2 Analysis of **TDescent**

In this section, we analyze TDescent, proving two main results. prove Proposition 5.6, which shows that TDescent terminates rapidly. Second, in Lemma 5.8 we show that the trust region constraint in Line 7 of Algorithm 3 (linesearch) prevents long steps. Thus, once the method enters a sufficiently small neighborhood of \overline{x} , it cannot leave.

We begin with descent Proposition 5.6, which relies on four technical lemmata that analyze the structure of Goldstein subgradients when σ is sufficiently small and x is sufficiently near \bar{x} : Lemma 5.2 states that elements of Goldstein subdifferential with small normal components are descent directions. Lemmas 5.3 and 5.4 show that normalized subgradient steps approximately reflect points across the active manifold. Lemma 5.5 uses the approximate reflection property to show that TDescent geometrically decreases the normal component of the input subgradient, ensuring that we rapidly find a descent direction. We now turn to the Lemmata.

5.2.1 Descent with Small Normal Part

The first lemma shows that Goldstein subgradients with small normal components are descent directions.

Lemma 5.2 (Descent with small normal part) *Define*

$$C_3 := \frac{C_1^2}{8L}.$$

Then for all $x \in B_{\delta_A}(\vec{x})$, $\sigma > 0$, and $g \in \partial_{\sigma} f(x) \setminus \{0\}$ satisfying

- 1. $\max\{\operatorname{dist}(x, M), \sigma\} \le \frac{C_2}{4} P_M(x) \overline{x};$ 2. $P_{NM}(P_M(x))(g) \le C_3 P_M(x) \overline{x}^2,$

we have

$$f \quad x - \sigma \frac{g}{g} \le f(x) - \frac{\sigma g}{8}$$

Proof We begin with preliminary notation and bounds. We fix $x \in B\delta_A(\vec{x})$ and subgradient $g \in \partial \sigma f(x) \setminus \{0\}$. We define g := PM(x), g := TM(y), and g := TM(y). We observe that

$$\sigma \leq \frac{C_2}{4} \ y - \bar{x} \leq \frac{C_2}{2} \ x - \bar{x} \leq C_2 \delta_{\mathrm{A}} \leq \delta_{\mathrm{A}},$$

where the final inequality follows since $C_2 \le 1$; see Lemma 4.2. We now turn to the proof.

The starting point of the proof is Lebourg's mean value Theorem [16, 2.4], which ensures that there exists $v \in \partial_{\sigma} f(x)$ such that

$$f \quad x - \sigma \frac{g}{g} \quad - f(x) = v, -\sigma \frac{g}{g} = -\frac{\sigma}{g} v, P_T(g) - \frac{\sigma}{g} v, P_N(g)$$

In what follows, we will show that the first term satisfies V, $P_T(g) \ge \frac{3}{8} g^2$, while the second term satisfies $|V, P_N(g)| \le \frac{1}{8} g^2$, yielding the result. Indeed, beginning with $|V, P_N(g)|$, we note that

$$P_N(g) \le C_3 P_M(x) - x^2 \le \frac{C_3}{C_1^2} g^2 = \frac{1}{8L} g^2,$$
 (5.1)

where the second inequality follow from Lemma 4.2. Consequently, we have the bound $|V, P_N(g)| \le L P_N(g) \le \frac{1}{8} g^{-2}$, where we first inequality relies on the estimate $V \le L$; see Item 6 of Proposition 3.5.

Next, we prove a lower bound on V, $P_T(g)$. Since $v \in \partial_{\sigma} f(x)$,

$$P_T(v-g) \leq \ 2C_{(a)}(\operatorname{dist}(x,M) + \sigma) \leq \ C_2C_{(a)} \ P_M(x) - \bar{x} \leq \ \frac{C_2C_{(a)}}{C_1} \ g \leq \ \frac{1}{2} \ g.$$

where the first inequality follows from (3.5); the second by assumption; the third follows from Lemma 4.2; and the fourth follows from the bound $\leq \frac{C_1}{2C(c)}$. Therefore,

$$P_T(v) - g \le P_T(v - g) + P_N(g) \le \frac{1}{2} g + \frac{1}{8L} g^2 \le \frac{5}{8} g$$

where the second inequality follows from (5.1) and the third follows from the bound $g \leq L$. Consequently, we have the bound

$$V, P_T(g) = P_T(V), g \ge g^2 - P_T(V) - g g \ge \frac{3}{8} g^2.$$

This completes the proof.

Note that the proof implies a slightly stronger bound than claimed, namely that we have $f(x - \sigma q/q) \le f(x) - \sigma q/4$. For the sake of maintaining symmetry with Proposition 5.1, however, we use the constant 1/8 throughout.

5.2.2 The Approximate Reflection Property

The next two lemmata prove the approximate reflection property that was described in the introduction. The lemmas roughly show that normalized subgradient steps approximately "flip the sign" of the normal component of the subgradient nearby the manifold; see Sect. 1.3.2 for more intuition. The first lemma proves the approximate reflection property up to a tolerance depending on the distance to the manifold and σ . This lemma will be used again in the proofs of Lemma 5.4 and Lemma 5.7.

Lemma 5.3 (Approximate reflection inequality, general case) For all $x \in B_{\delta_A/2}(x)$, $\sigma \in (0, \delta_A/2]$, $g \in \partial_\sigma f(x) \setminus \{0\}$ and $\hat{g} \in \partial_\sigma f(x) \setminus \{0\}$, we have

$$P_{NM}(P_{M}(x))(\hat{g}), g$$

$$\leq -\mu \quad P_{NM}(P_{M}(x))g + \frac{(\mu + L) \ g \ \operatorname{dist}(x, M)}{\sigma}$$

$$+ \frac{(\mu + L) \ g \ CM \ (\operatorname{dist}^{2}(x, M) + \sigma^{2})}{\sigma}. \tag{5.2}$$

Proof We begin with preliminary notation and bounds. We fix $x \in B_{\delta_A/2}(\overline{x})$ and subgradient $g \in \partial \sigma f(x) \setminus \{0\}$. We define y := PM(x), T := TM(y), and N := NM(y). Finally, define $u := \frac{g}{g}$. Note that since $x \in B_{\delta_A/2}(\overline{x})$ and $\sigma \le \delta_A/2$, we have $x - \sigma u \in B_{\delta_A}(\overline{x})$.

Therefore, by the aiming inequality (3.6), we have

$$\|\hat{g}_{x} \times -\sigma u - PM (x - \sigma u)\|_{\infty}^{2} \leq \mu \underbrace{x - \sigma u - PM (x - \sigma u)}_{=: B} \%$$

We aim to simplify this inequality with (3.1). To that end, first note that

$$(x - \sigma u - P_M (x - \sigma u)) - (x - P_M (x) - \sigma P_N (u))$$

$$= P_M (x - \sigma u) - P_M (x) + \sigma P_T (u)$$

$$\leq C_M (\text{dist}^2(x, M) + \sigma^2).$$
(5.3)

Consequently, we have

$$A \geq B \geq \mu \quad x - P_M(x) - \sigma P_N(u) - \mu \quad C_M(\operatorname{dist}^2(x, M) + \sigma^2) \geq \sigma \mu \quad P_N(u) - \mu \quad S_n(u) = 0$$

where $S := \operatorname{dist}(x, M) + CM \left(\operatorname{dist}^2(x, M) + \sigma^2\right)$. In addition, by (5.3) we have

$$\hat{g}$$
, $(x - \sigma u - P_M (x - \sigma u)) + \sigma P_N u \le L S$.

Therefore, we have

$$\hat{g}, \sigma P_N(u) = -A + \hat{g}, (x - \sigma u - P_M(x - \sigma u)) + \sigma P_N u$$

$$\downarrow \hat{\Box} \hat{\Box} \hat{\Box}$$

$$\leq -\sigma\mu \quad P_N(u) + (\mu + L)S. \tag{5.4}$$

Inequality (5.2) then follows by multiplying both sides of inequality (5.4) by g/σ .

The second lemma is an application of Lemma 5.3 nearby the manifold.

Lemma 5.4 (Approximate reflection inequality near the manifold) Define

$$C_4 := \min \ \frac{\beta}{C_{(a)}(1+\delta_{\rm A})}, \frac{\min\{\mu/\delta_{\rm A}, \, C_3D_2/\beta\}}{4(1+(1+\delta_{\rm A})C_M)(\mu+L)}, \frac{1}{2} \ .$$

Then for all $x \in B\delta_A/2(x)$, $\sigma > 0$, and $g \in \partial \sigma f(x) \setminus \{0\}$ satisfying

$$\max \ \frac{\operatorname{dist}(x, M)}{\sigma}, \sigma \le C_4 \ P_M(x) - x,$$

we have

$$\begin{split} P_{N_M \ (P_M \ (x))}(\hat{g}), \ g &\leq - \ D_2 \ P_{N_M \ (P_M \ (x))}g + \frac{C_3D_2}{2} \ P_M \ (x) - \frac{1}{x}^{-2} \\ & \text{for all } \hat{g} \in \partial \ f \quad x - \sigma \frac{g}{g} \quad . \end{split}$$

Proof We begin with preliminary notation and bounds. We fix $x \in B_{\delta_A/2}(\overline{x})$ and subgradient $g \in \partial \sigma f(x) \setminus \{0\}$. We define $y := P_M(x)$, $T := T_M(y)$, and $N := N_M(y)$. We observe that

$$\sigma \le C_4 \quad y - \overline{x} \le 2C_4 \quad x - \overline{x} \le C_4 \delta_A \le \delta_A / 2$$

Finally, we have

$$S := \operatorname{dist}(x, M) + C_M \left(\operatorname{dist}^2(x, M) + \sigma^2\right) \le \sigma C_4 (1 + C_M (1 + \delta_A)) \quad y - \overline{x}.$$

$$(5.5)$$

where the inequality follows from the bound dist $(x, M) \le x - \overline{x} \le \delta$ A. We now apply inequality (5.2):

$$\begin{split} P_N \hat{g}, \, g & \leq -\mu \quad P_N g + \ \frac{(\mu + \ L) \ g \ S}{\sigma} \\ & \leq -\mu \quad P_N g + (\ 1 + (1 + \delta_A) C_M \) (\mu + \ L) C_4 \ g \ y - \overline{x} \\ & \leq -\mu \quad P_N g + (\ 1 + (1 + \delta_A) C_M \) (\mu + \ L) C_4 (\ P_T (g) + \ P_N (g)) \ y - \overline{x} \\ & \leq -\frac{\mu}{2} \ P_N g + \ \frac{C_3 D_2}{4\beta} \ P_T (g) \ y - \overline{x}, \end{split}$$

where the second inequality follows from (5.5); the third inequality follows from triangle inequality; and the fourth inequality follows from the bound

$$(1 + (1 + \delta_{A})C_{M})(\mu + L)C_{4} y - \bar{\chi} \le \frac{\mu/\delta_{A}}{4} \cdot (2 x - \bar{\chi}) \le \frac{\mu/\delta_{A}}{4} \delta_{A} \le \mu/2.$$

The proof will be complete if we can show that

$$P_T(g) \le 2\beta y - x$$
.

To that end, we have

$$P_T(g) \le C_{(a)}(\operatorname{dist}(x, M) + \sigma) + \beta \quad y - \overline{x}$$

$$\le (C_4C_{(a)}(1 + \delta_A) + \beta) \quad y - \overline{x} \le 2\beta \quad y - \overline{x},$$

where the first inequality follows from (3.4); the second inequality follows from the lemma assumptions and the bound dist $(x, M) \le C_4 \sigma y - \bar{x} \le C_4 \delta_A y - \bar{x}$; and the third inequality follows from the bounds on C_4 . This completes the proof.

5.2.3 The Normal Component Shrinks Geometrically

The following lemma shows that every step of TDescent geometrically shrinks the normal component of the subgradient, up to a tolerance of $O(PM(x) - x^2)$.

Lemma 5.5 (Normal component shrinks geometrically) Define

$$C_5 := \min \ \frac{\beta}{2C_{(a)}}, \frac{C_3D_2}{32C_{(a)}\beta}, C_4, \frac{C_2}{4}$$
.

Then for all $x \in B_{\delta_A/2}(x)$, $\sigma > 0$, $g \in \partial \sigma f(x) \setminus \{0\}$, and $\hat{g} \in \partial f(x - \sigma - \frac{g}{g}) \setminus \{0\}$ satisfying

1.
$$P_{NM}(P_{M}(x))g \ge C_3 P_{M}(x) - \overline{x}^2$$
;

1.
$$P_{NM (P_M (x))} g \ge C_3 P_M (x) - \overline{x}^2$$
;
2. $\max \frac{\operatorname{dist}(x, M)}{\sigma}$, $\sigma \le C_5 P_M (x) - \overline{x}$,

the vector $g = \operatorname{argmin}_{h \in [q, \hat{q}]} h$ satisfies:

$$P_{N_M (P_M (x))}(g)^2 \le 1 - \frac{3D_2^2}{64L^2} P_{N_M (P_M (x))}g^2.$$

Proof We begin with preliminary notation and bounds. We fix $x \in B_{\delta_A}/2(x)$ and subgradient $g \in \partial \sigma f(x) \setminus \{0\}$. We define y := PM(x), T := TM(y), and N := TM(y)NM (y). We observe two bounds. First, we have

$$\sigma \leq C_5 \ y - \bar{x} \leq \ 2C_5 \ x - \bar{x} \leq \ C_5 \delta_{\mathrm{A}} \leq 1.$$

where the final inequality follows since $C_5 \le C_2/4 \le 1/(8\delta_A)$. Second, we have

$$\operatorname{dist}(x, M) \le C_5 \sigma \ y - \overline{x} \le C_5 \ y - \overline{x}, \tag{5.6}$$

since $\sigma \leq 1$. We now turn to the proof.

Consider the optimal weight λ := $\operatorname{argmin}_{\lambda \in [0,1]} g + \lambda (\hat{g} - g)$. By definition we have $g = g + \lambda (\hat{g} - g)$. Moreover, a quick calculation shows that

$$\lambda = \max \min -\frac{g, \hat{g} - g}{\hat{q} - q^2}, 1 , 0 .$$

We claim that the following bound holds on λ :

$$-\frac{P_{N}(g), \hat{g} - g}{\|\hat{g}\|_{2}^{2}} \le \lambda \le -\frac{3 P_{N}(g), \hat{g} - g}{\|2 P_{N}(\hat{g}\|_{2}^{2} - g)^{2}}.$$

$$=:\lambda_{1}$$

$$=:\lambda_{2}$$
(5.7)

Note that (5.7) is an immediate consequence of the following bound:

$$0 \le -\frac{1}{2} P_N(g), \, \hat{g} - g \le -g, \, \hat{g} - g \le -\frac{3}{2} P_N(g), \, \hat{g} - g \,. \tag{5.8}$$

Indeed, if (5.8) holds, then $\lambda = \min - \frac{g_1 \hat{g} - g_2}{\hat{g} - g_2}$, 1 · Thus, we obtain the upper bound

$$\lambda \leq -\frac{g, \hat{g} - g}{\hat{g} - g^{2}} \leq -\frac{3}{2} \frac{P_{N}(g), \hat{g} - g}{g - \hat{g}^{2}} \leq -\frac{3}{2} \frac{P_{N}(g), \hat{g} - g}{P_{N}(g - \hat{g})^{2}} = \lambda_{2}.$$

Likewise, we obtain the lower bound

$$\lambda = \min -\frac{g, \hat{g} - g}{\hat{g} - g^{2}}, 1 \ge \min -\frac{g, \hat{g} - g}{4L^{2}}, 1$$

$$= -\frac{g, \hat{g} - g}{4L^{2}} \ge -\frac{P_{N}(g), \hat{g} - g}{8L^{2}} = \lambda_{1},$$

where the first inequality follows from the bound $\hat{g} - g^2 \le 2(\hat{g}^2 + g^2) \le 4L^2$; and the second equality follows from the bound $|g, \hat{g} - g| \le g^2 = 2L^2$. Thus, we now prove (5.8).

To that end, note that (5.8) is equivalent to the following bound:

$$P_T(g), \hat{g} - g \leq \frac{-P_N(g), \hat{g} - g}{2}.$$
 (5.9)

Therefore, we first bound ${}^{\&}_{P_T(g)}$, $\hat{g} - g {}^{\&}_{P_T(g)}$

$$P_T(g), \hat{g} - g \leq P_T(g) \quad P_T(\hat{g} - g)$$

$$\leq 2C_{(a)}(\operatorname{dist}(x, M) + \sigma)(C_{(a)}(\operatorname{dist}(x, M) + \sigma) + \beta \quad y - \overline{x})$$

$$\leq 4C_{(a)}C_{5}(2C_{(a)}C_{5} + \beta) \quad y - \overline{x}^{2}$$

$$\leq \frac{C_{3}D_{2}}{4} \quad y - \overline{x}^{2},$$

where the second inequality follows from (3.4) and (3.5); the third inequality follows from (5.6) and the bound $\sigma \le C_5 \ y - \overline{x}$; and the fourth inequality follows from the definition of C_5 . To complete the proof of (5.9), we show that $\frac{C_3D_2}{4} \ y - \overline{x}^2 \le -\frac{1}{2} \ P_N(g)$, $\hat{g} - g$:

$$\frac{C_3 D_2}{2} y - \bar{x}^2 \le D_2 P_N(g) - \frac{C_3 D_2}{2} y - \bar{x}^2$$

$$\le - P_N(g), \hat{g}$$

$$\le - P_N(g), \hat{g} - g, \qquad (5.10)$$

where the first inequality follows from the assumption $\frac{D_2}{2}$ $P_N(g) \ge \frac{C_3D_2}{2}$ $y - \overline{x}^2$; the second inequality follows from Lemma 5.4 (recall $C \le C_4$ and $x \in B_{\delta_A/2}(x)$); and the third inequality follows from $P_N(g)$, $g = P_N(g)^2 \ge 0$. Thus, the equivalent bounds (5.9) and (5.8) hold. Consequently, Equation (5.7) holds.

Now we turn to the contraction argument. Consider the function $r: R \to R$ satisfying

$$r(\lambda) = P_N(g)^2 + 2\lambda P_N(g), \hat{g} - g + \lambda^2 P_N(\hat{g} - g)^2$$
 for all $\lambda \in \mathbb{R}$.

Observe that

$$P_N(g)^2 = P_N(g)^2 + 2\lambda P_N(g), \hat{g} - g + (\lambda)^2 P_N(\hat{g} - g)^2 = r(\lambda).$$

Therefore, by convexity of r and (5.7), we have

$$P_N(g)^2 = r(\lambda) \le \max_{\lambda \in [\lambda_1, \lambda_2]} r(\lambda) \le \max\{r(\lambda_1), r(\lambda_2)\}.$$

To complete the proof, we show each term in the "max" is bounded by $1 - \frac{3D_2^2}{64L^2} - P_N(g)^2$.

To show this, we will use the following consequence of (5.10):

$$- P_N(g), \hat{g} - g \ge D_2 P_N(g) - \frac{C_3 D_2}{2} y - \bar{x}^2 \ge \frac{D_2}{2} P_N(g), \tag{5.11}$$

where the final inequality follows from the assumption $\frac{C_3D_2}{2}$ $y-\bar{x}^2 \leq \frac{D_2}{2}$ $P_N(g)$. Indeed, first observe that

$$r(\lambda_2) = P_N(g)^2 - \frac{3}{4} \frac{P_N(g), \hat{g} - g^2}{P_N(\hat{g} - g)^2}$$

$$\leq 1 - \frac{3D_2^2}{16 P_N(\hat{g} - g)^2} P_N(g)^2$$

$$\leq 1 - \frac{3D_2^2}{64L^2} P_N(g)^2,$$

where the first inequality from (5.11) and the second inequality follows from the bound $P_N(\hat{g} - g)^2 \le \hat{g} - g^2 \le 4L^2$. Likewise, observe that

$$\begin{split} r\left(\lambda_{1}\right) &= & P_{N}(g)^{2} - \frac{P_{N}(g), \, \hat{g} - g^{2}}{4L^{2}} + \frac{P_{N}(g), \, \hat{g} - g^{2}}{64L^{4}} P_{N}(\hat{g} - g)^{2}}{64L^{4}} \\ &\leq & P_{N}(g)^{2} - \frac{P_{N}(g), \, \hat{g} - g^{2}}{4L^{2}} + \frac{P_{N}(g), \, \hat{g} - g^{2}}{16L^{2}} \\ &\leq & 1 - \frac{3D_{2}^{2}}{64L^{2}} P_{N}(g)^{2}, \end{split}$$

where the first inequality follows from the bound $P_N(\hat{g} - g)^2 \le \hat{g} - g^2 \le 4L^2$; and the second inequality follows from (5.11). Therefore, the proof is complete.

5.2.4 TDescent Terminates with Descent

The following proposition is the main result of this section. It shows that TDescent must either terminate with descent or f(x) - f(x) is already exponentially small in T.

Proposition 5.6 (TDescent loop terminates with descent) $Fix T \in \mathbb{N}$. Then for all $x \in B\delta_A/2(x)$, $v \in \partial\sigma f(x)$, and $\sigma > 0$ satisfying

$$\max \ \frac{\operatorname{dist}(x, M)}{\sigma}, \sigma \le C_5 \ P_M(x) - \overline{x},$$

at least one of the following holds:

1. we have

$$f(x) - f(\bar{x}) \le \frac{(C_5^2 L + \beta) L}{C_3} \quad 1 - \frac{3\mu^2}{256L^2} \quad ;$$

2. the vector $g := TDescent(x, V, \sigma, T)$ satisfies g > 0 and

$$f \quad x - \sigma \frac{g}{g} \le f(x) - \frac{\sigma \operatorname{dist}(0, \, \partial_{\sigma} f(x))}{8}.$$

Proof We begin with preliminary notation and bounds. We fix $x \in B\delta_A/2(x)$ and subgradient $v \in \partial \sigma f(x)$. We define y := PM(x), and N := NM(y). Observe that

$$\sigma \le C_5 \ y - \overline{x} \le 2C_5 \ x - \overline{x} \le C_5 \delta_A \le 1$$

where the final inequality follows by definition of $G \le C_2/4 \le 1/(8\delta_A)$. In addition, since $C_5 \le C_2/4$, we have $\sigma \le (C_2/4)$ $y = \overline{x}$ and

$$\operatorname{dist}(x,\,M) \leq \sigma \ C_5 \ y - x \leq \frac{C_2}{4} \ y - x,$$

where the final inequality follows from $\sigma \leq 1$. Consequently,

$$\max\{\sigma, \operatorname{dist}(x, M)\} \le \frac{C_2}{4} y - \overline{x}. \tag{5.12}$$

We now turn to the proof.

Turning to the proof, note that since $x \in B\delta_A/2(x)$, Lemma 4.2 and (5.12) ensure that

$$\operatorname{dist}(0,\,\partial_\sigma\,f(x))\geq \,C_1\,\,y-\overline{x}>\,\,0\cdot$$

Thus, if $\mathsf{TDescent}(x, V, \sigma, T)$ terminates at $t \leq T$, then Item 2 must hold. For the remainder of the proof, we suppose that $\mathsf{TDescent}(x, V, \sigma, T)$ terminates at the final iteration t = T and that Item 2 does not hold. In this case, Lemma 5.2 and (5.12) ensure that the iterates g_t of $\mathsf{TDescent}(x, V, \sigma, T)$ satisfy $P_N(g_t) > C_3 \ y - \overline{x}^2$ for all $0 \leq t \leq T - 1$. Therefore, since $x \in B\delta_A/2(\overline{x})$, $\max\{\mathrm{dist}(x, M)/\sigma, \sigma\} \leq C_5 \ y - \overline{x}$, and $P_N(g_t) > C_3 \ y - \overline{x}^2$, Lemma 5.5, yields the contraction:

$$P_N(g_{t+1})^{-2} \le 1 - \frac{3D_2^2}{64L^2} - P_N(g_t)^{-2}$$
, for all $0 \le t \le T - 1$.

Unfolding this contraction, we see that g_T is an exponentially small Goldstein subgradient:

$$P_N(g_T) \le 1 - \frac{3D_2^2}{64L^2} P_N(g_0).$$

As a result, the projection *y* is nearby \bar{x} :

$$y - \bar{\chi}^2 \le \frac{P_N(g_T)}{C_3} \le \frac{P_N(g_0)}{C_3} \quad 1 - \frac{3D_2^2}{64L^2} \quad ^{T/2} \le \frac{L}{C_3} \quad 1 - \frac{3D_2^2}{64L^2} \quad .$$
 (5.13)

Consequently,

$$\begin{split} f(x) - f(\bar{x}) &\leq L \mathrm{dist}(x, M) + \beta \quad y - \bar{x}^{-2} \\ &\leq (C_5^2 L + \beta) \quad y - \bar{x}^{-2} \\ &\leq \frac{(C_5^2 L + \beta) L}{C_3} \quad 1 - \frac{3D_2^2}{64L^2} \end{split},$$

where the first inequality follows from (3.7) (recall $\not = B\delta_A/2(\vec{x})$); the second inequality follows since dist(x, M) $\leq \sigma$ C_5 $y - \bar{x} \leq C_5^2$ $y - \bar{x}^2$; and the third inequality follows from (5.13). The proof then follows from the identity $D_2 = \frac{\mu}{2}$.

5.2.5 The "Trust Region" Constraint Prevents Long Steps

Before ending this section, we must establish one final technical result for Descent. Namely, in Lemma 5.8, we show that for appropriate, TDescent eventually generates small subgradients on the order of O(x - x). This property is intuitive because $\operatorname{dist}(0, \partial_{\sigma} f(x)) = 0$ whenever $\sigma \ge x - x$. This property will help us ensure that the iterates of NTDescent (Algorithm 4) cannot leave sufficiently small neighborhoods of x. Indeed, since the subgradients V_{i+1} generated by Algorithm 3 (linesearch) are decreasing in norm, we will show that the trust region constraint $\sigma_i \le \frac{V_{i+1}}{s}$ in Line 7 of Algorithm 3 must eventually be violated for large i. This ensures large σ_i are never chosen.

To prove this claim, we first establish a refinement of the approximate reflection property in Lemma 5.4. Compared to Lemma 5.4, the following lemma deals with a different range of parameters. We place the proof in Appendix D as it follows from a similar line of reasoning as Lemma 5.4.

Lemma 5.7 (Approximate reflection across manifold, large steps) *Define*

$$\delta_{\text{Grid}} := \min \frac{\delta_{\text{A}}}{2}, \frac{1}{C_M (D_1^{-1} + 1)}, \frac{\mu}{8(C_{(a)} + \beta)}$$
.

Then for all $x \in B\delta_{Grid}(x)$, $\sigma > 0$, and $g \in \partial_{\sigma} f(x) \setminus \{0\}$ satisfying

$$D_1^{-1} \operatorname{dist}(x, M) \le \sigma \le \delta$$
 Grid

we have

$$\hat{g}, g \leq -D_2 \ g + 2D_2 \ P_{T_M} \left(P_M \left(x \right) \right) \left(g \right) \quad \text{ for all } \hat{g} \in \partial \ f \quad x - \sigma \frac{g}{g} \quad .$$

Finally, we prove that TDescent eventually generates small subgradients.

Lemma 5.8 (TDescent yields small subgradients) $Fix\ T\in \mathbb{N}$. Then for all $x\in B\delta_{Grid}(\vec{x}), \sigma \geq 0$, and $g\in \partial\sigma\ f(x)\setminus\{0\}$ satisfying

$$D_1^{-1} \operatorname{dist}(x, M) \le \sigma \le \delta$$
 Grid

the vector $g := TDescent(x, g, \sigma, T)$ satisfies

$$g \leq \max \left\{ \begin{array}{ccc} & \mu^2 & T'^2 \\ 1 - \frac{\mu^2}{64L^2} & g, & 4C_{(a)}\sigma + 4(C_{(a)} + 2\beta) & x - \overline{x}, & \frac{8(f(x) - f(\overline{x}))}{\sigma} \end{array} \right\}.$$

Proof We begin with preliminary notation and bounds. We fix $x \in B\delta_{Grid}(\overline{x})$ and subgradient $g \in \partial \sigma f(x)\setminus\{0\}$. We define y := PM(x) and T := TM(y). We also define $c := C(a)(\operatorname{dist}(x, M) + \sigma) + \beta \quad y - \overline{x}$. We have the following two bounds: First, we have

$$c \le C_{(a)}(x - x + \sigma) + 2\beta x - x \le C_{(a)}\sigma + (C_{(a)} + 2\beta) x - x.$$
 (5.14)

Second, by (3.4), we have

$$P_T(v) \le c \quad \text{for all } v \in \partial_{\sigma} f(x).$$
 (5.15)

We now turn to the proof.

Note that the result holds automatically if g = 0. Thus, we first consider the case where TDescent terminates in descent, meaning

$$f(x_+) - f(x) \le -\frac{\sigma g}{8}$$
 where $x_+ := x - \sigma \frac{g}{g}$.

Since $\sigma \le \delta_{\text{Grid}} \le \delta_{\text{A}}/2$ and $x \in B_{\delta_{\text{A}}/2}(\vec{x})$, it follows that $x + \in B_{\delta_{\text{A}}}(\vec{x})$. Thus, by Item 1 of Proposition 3.5, we have

$$f(x_+) \ge f(\bar{x}) + \frac{y}{2} |x - \bar{x}|^2 \ge f(\bar{x}).$$

Consequently, we have

$$f(\bar{x}) - f(x) \le -\frac{\sigma g}{8}$$

Rearranging then gives the upper bound $g \le \frac{8(f(x) - f(\bar{x}))}{\sigma}$, as desired.

Let us now suppose that TDescent does not terminate with descent or with g = 0. In this case, the iterates g_0, \dots, g_T of TDescent (x, g, σ, T) exist and satisfy $g_t \in \partial \sigma f(x)$ for all $t \leq T$. We consider two cases.

Case 1. Now suppose $g_t \le 4c$ for some t satisfying $0 \le t \le T$. Since g_t is a decreasing sequence, it follows that $g = g_T \le 4c$. Recalling (5.14), yields the bound

$$g \le 4c \le 4C_{(a)}\sigma + 4(C_{(a)} + 2\beta) x - x$$

as desired.

Case 2. Next suppose that for all $0 \le t \le T$ we have $4c \le g_t$. In this case, Lemma 5.7 shows that for all $t \le T$, we have

$$\hat{g}_t, g_t \le -\frac{\mu}{2} g_t + \mu \quad P_T g_t \le -\frac{\mu}{2} g_t + \mu \ c \le -\frac{\mu}{4} g_t. \tag{5.16}$$

We now use this bound to prove a one-step geometric improvement bound for g_t^2 . To that end, fix any $t \le T - 1$ and define the weight $\lambda := \frac{\mu g_t}{16L^2}$ and the vector $g_{\lambda} := g_t + \lambda(\hat{g}_t - g_t)$. Notice that $\lambda \in [0, 1]$, since

$$\lambda = \frac{\mu \ g_t}{16L^2} \le \frac{\mu}{16L} \le 1,$$

where the first equation follows since $g_t \in \partial \sigma f(x)$ and the second follows since $L \ge \mu$; see Lemma E.1. Thus

$$\begin{split} g_{t+1} \ ^2 \leq \ & g_{\lambda} \ ^2 = \ g_t \ ^2 + 2\lambda \ g_t, \hat{g}_t - g_t + \lambda^2 \hat{\ } g_t - g_t \ ^2 \\ \leq \ & g_t \ ^2 + 2\lambda \ g_t, \hat{g}_t - 2\lambda \ g_t \ ^2 + 4L^2\lambda^2 \\ \leq \ & g_t \ ^2 - \frac{\lambda\mu}{2} \ g_t + 4L^2\lambda^2 \\ = \ & 1 - \frac{\mu^2}{64L^2} \ g_t \ ^2, \end{split}$$

where the first inequality follows by definition of g_{t+1} ; the second inequality follows from the fact that L is a local Lipschitz constant of f near \overline{x} ; and the third inequality follows from (5.16). Thus, to complete the proof, simply unfold this recursion to get the bound

$$g = g_T \le 1 - \frac{\mu^2}{64L^2} g_0^{-2},$$

as desired.

6 Rapid Local Convergented escent

In this Section, we present our main convergence guarantees for the NTDescent method under Assumption A. The main results of the section are Theorem 6.3 and Theorem 6.5, which analyze the nonconvex and convex settings respectively. In the

Parameter	Definition
s _{lb}	c_0 g_0
a_1	$\min\{D_1, D_2/L\}$
a_2	$\frac{\min\{C_1/L,C_5\}}{2}$
δ_{LS}	$\min \frac{\delta_{A}}{2}, \ \delta_{GI}, \ \delta_{ND}, \ \delta_{Grid}, \ \frac{1}{2(a_{1}+2a_{2})}, \ \frac{\textit{V} \ \textit{D}_{1}^{2} \min\{\delta_{Grid}/2,1/4\}^{2}}{2L}, \ 1$
C_6	$\max 1, \frac{8(C_{(a)} + 2\beta + 2C_{(a)}D_1^{-1})}{s_{ b}}, 2D_1^{-1}, \frac{4\gamma D_1}{s_{ b}}$
1, T	$\max \ \frac{(C_5^2L+\beta)\ L}{C_3} \ 1 - \frac{3\mu^2}{256L^2} \ , \ 1 - \frac{\mu^2}{64L^2} \ L$
2, G	$\max \frac{L}{\min\{1,a_1\}} + \frac{\beta}{2\min\{1,a_1\}a_2^2}, 8C(a), L 2^{-G}$
ρ	$1 - \frac{1}{8} \min \frac{y_{a_2}}{8 \max\{4La_2^2, \beta\}}, \frac{\mu_{a_1}}{4 \max\{2L, \beta', a_2^2\}}$

Table 2 Parameters used throughout Sect. 6; see also Table 1

nonconvex setting, we prove that iterates of NTDescent locally nearly linearly converge, provided some iterate reaches a sufficiently small neighborhood of \overline{x} . In the convex setting, we strengthen this guarantee, showing that for any initial starting point x_0 and any failure probability p, there exists some index K_p after which NTDescent nearly converges linearly with probability at least f f Both results are a consequence of the local one-step improvement bound of Proposition 6.1. This proposition shows that with high probability, the following hold locally for linesearch: its output is nearby its input; and the function gap geometrically decreases whenever it is larger than a quantity that is exponentially small in the inner loop budget and the grid size. The former property will help ensure that the iterates of NTDescent do not escape a local neighborhood of \overline{x} .

6.1 Assumptions and Notation

Throughout this section, we assume the following assumptions and notations are in force. We assume that

- 1. the budget T_k and grid size G_k satisfy $\min\{T_k, G_k\} \ge k+1$ for all $k \ge 0$.
- 2. We fix an initial we an initial point $x_0 \in \mathbb{R}^d$ and $g_0 \in \partial f(x_0)$. We assume that $g_0 = 0$. We assume Assumption A is in force at a point $\overline{x} \in \mathbb{R}^d$ and use the notation of Proposition 3.5 throughout. We let $\{x_k\}$ denote the sequence of iterates generated by NTDescent $\{x_0, g_0, c_0, \{G_k\}, \{T_k\}\}$ when applied to f.

Turning to notation, we now summarize in Table 2 the main constants used in this section.

In the following, we lower and upper bound the trust region parameter in linesearch:

$$s_{lb} \le \max\{ g_k, c_0 g_0 \} \le L,$$
 (6.1)

where the lower bound follows by definition, and the upper bound follows from Part 6 of Proposition 3.5. In addition, we apply Theorem 4.3 with the constants $\boldsymbol{\varrho}$ a_2 . These constants are derived from the parameters D_1 , D_2 , C_1 , and C_5 which are defined in Lemmas 4.1, 4.2, and 5.5 respectively. We also define a neighborhood B $\delta_{LS}(\vec{x})$ for which linesearch results in geometric improvement. Here, the radius δ_{LS} is derived from the parameters δ_A , δ_{GI} , δ_{ND} , δ_{Grid} , and γ which appear in Proposition 3.5 and Lemmas 4.1, 5.1, 5.7, and 5.8. In addition, the constant G will appear in an upper bound on the steplength of linesearch.

We then define three terms $_{1,T}$, $_{2,G}$, and ρ which appear in our convergence rate analysis. These terms are defined for all $_{T}$, $_{G}$ > 0 and are derived from the parameters C_5 , C_3 , a_1 , a_2 , L, β , $C_{(a)}$, γ , and μ which appear in Lemma 5.2, Lemma 5.5, Proposition 3.5, and Assumption A.

Finally, in the following propositions, the constant $\rho \in (0, 1)$ plays the role of a local contraction factor, while the terms $_{1,T}$ and $_{2,G}$ are upper bounds for function gap of NTDescent.

We now turn to the one-step improvement argument.

6.2 One Step Improvement

The following proposition presents our one-step improvement bound.

Proposition 6.1 (One step improvement) Assume the assumptions of Sect. 6.1 are satisfied. Recall the notation in Table 2. Then the following holds for all $x \in B\delta_{LS}(\vec{x})$, subgradients $g \in \partial f(x)$, and grid sizes $G > \log_2(1/\delta_{Grid})$: Fix a scalar $s \in [s_{lb}, L]$, a failure probability $p \in (0, 1)$ and budget T satisfying

$$T \ge \frac{256L^2}{\mu^2} - 2\log(1/p)$$
.

Then with probability at least 1 - p, the point $\tilde{x} = linesearch(x, g, s, G, T)$ satisfies

1.
$$f(\tilde{x}) - f(\bar{x}) \le \max\{\rho(f(x) - f(\bar{x})), \sqrt{1,T'-2,G}\};$$

2. $\tilde{x} - x \le C_6 \max_{1,T} \sqrt{s_{lb}}, 2, G's_{lb}, \frac{2(f(x) - f(\bar{x}))}{2(f(x) - f(\bar{x}))}$,

Proof We fix $x \in B_{\delta_{LS}}(x)$, define y := PM(x), and choose a subgradient $g \in \partial f(x)$. Throughout we may freely use the results of Proposition 3.5 since $\delta_{LS} \le \delta_A$. We will first establish the first item of the Proposition. To that end, let us assume that

$$f(x) - f(\bar{x}) > \max\{_{1,T},_{2,G}\};$$

otherwise the proof is trivial. In this case, we claim that x must satisfy either Item 1 or Item 2 of Theorem 4.3 for at least one G_i with $i \le G - 1$. To derive a contradiction, suppose that both items are not satisfied for x with any choice G_i with $i = 0, \dots, G - 1$. We will show that neither Item 1b nor its complement can be satisfied, leading to a contradiction.

Throughout the following argument, we will use the following bound:

$$\max\{a_1 \operatorname{dist}(x, M), a_2 \ y - \overline{x}\} \le (a_1 + 2a_2) \delta_{LS} \le \frac{1}{2} = \sigma_{G-1}.$$

Now suppose that Item 1b holds, i.e., $\frac{2}{4} y - \overline{x}^2 \le \operatorname{dist}(x, M)$. Then by assumption, Item 1a must fail for $\operatorname{any} \sigma_i$. We claim that this failure ensures that $\sigma_0 \ge a_1 \operatorname{dist}(x, M)$. Indeed, if $\sigma_0 \le a_1 \operatorname{dist}(x, M)$, we must have

$$\sigma_0 \le (a_1/2) \operatorname{dist}(x, M) \le a_1 \operatorname{dist}(x, M) \le \sigma_{G-1}$$

since σ_0 cannot satisfy Item 1a. Thus, there exists some $j \le G - 1$ such that $\sigma_j = 2^j \sigma_0$ satisfies Item 1a, a contradiction. Therefore, we have

$$\sigma_0 > a_1 \text{dist}(x, M) \ge a_1 a_2^2 y - \bar{x}^2$$

In this case, by (3.7), we have

$$f(x) - f(\bar{x}) \le L \operatorname{dist}(x, M) + \frac{\beta}{2} |y - \bar{x}|^2 \le \frac{L}{a_1} + \frac{\beta}{2a_2^2 a_1} |\sigma_0 \le \alpha_2 G$$

which is a contradiction. Therefore, Item 1b cannot hold, so we have $a_2^2 y - \bar{x}^2 > \text{dist}(x, M)$.

Next, for the sake of contradiction, suppose that there exists σ_i satisfying Item 2a. In this case, since $\sigma_i \ge (a_2/2)$ $y - \bar{x}$, we have

$$dist(x, M) < a_2^2 y - \bar{x}^2 \le 2a_2\sigma_i y - \bar{x}$$

i.e., σ_i also satisfies Item 2b, which is a contradiction. Therefore no σ_i satisfies Item 2a. We claim that this ensures $\sigma_0 > a_2 \ y - \overline{x}$. Indeed, if $\sigma_0 \le a_2 \ y - \overline{x}$, we must have

$$\sigma_0 \le (a_2/2) \ y - \overline{x} \le a_0 \ y - \overline{x} \le \sigma_{G-1}$$

since σ_0 cannot satisfy Item 2a. Thus, there exists some $j \le G - 1$ such that $\sigma_j = 2^j \sigma_0$ satisfies Item 2a, a contradiction. Therefore, we have

$$\sigma_0 > a_2 \ y - \bar{x} \ge \overline{\operatorname{dist}(x, M)}$$
.

In this case, by (3.7), we have

$$f(x) - f(\bar{x}) \le L \text{dist}(x, M) + \frac{\beta}{2} y - \bar{x}^2 \le L + \frac{\beta}{2a_2^2} \sigma_0^2 \le 2.6$$

which is a contradiction. Therefore, there must exist \mathcal{O}_i satisfying either Item 1 or Item 2 of Theorem 4.3.

Let us now fix a σ_i satisfying either Item 1 or Item 2 of Theorem 4.3. Then, by Theorem 4.3, we have the bound

$$\sigma_i \operatorname{dist}(0, \partial_{\sigma_i} f(x)) \ge 8(1 - \rho)(f(x) - f(\overline{x})).$$

$$0 < \sigma_i \le a_1 \operatorname{dist}(x, M) \le D_1 \operatorname{dist}(x, M). \tag{6.2}$$

Finally, from the definition $D_2 = \mu/2$, it follows that T satisfies the conditions of Proposition 5.1. Therefore, with probability at least 1 - p, we have

$$f - x - \sigma_i \frac{V_{i+1}}{V_{i+1}} - f(\bar{x}) \le f(x) - f(\bar{x}) - \frac{\sigma_i}{8} \mathrm{dist}(0, \, \partial_{\sigma_i} f(x)) \le \rho(f(x) - f(\bar{x})).$$

Next, we show that V_{i+1} and σ_i satisfy the trust region condition $\sigma_i \leq \frac{V_{i+1}}{s}$. To that end, note that the conditions of Lemma 4.1 are met: We have $x \in B\delta_{GI}(x)$ since $\delta_{LS} \leq \delta_{GI}$. We also have bound $\sigma_i \leq D_1 \mathrm{dist}(x, M)$ from (6.2). Therefore, it follows that the minimal norm Goldstein subgradient is lower bounded: $\mathrm{di}(\Omega, \partial_{\sigma_i} f(x)) \geq D_2$. Consequently, we have

$$\sigma_i \leq a_1 \mathrm{dist}(x,\,M) \leq \frac{D_2 \delta_{\mathsf{LS}}}{s} \leq \frac{\mathrm{dist}(0,\,\partial_{\sigma_i} \, f(x)) \delta_{\mathsf{LS}}}{s} \leq \frac{v_{i+1}}{s},$$

where the second inequality follows from the definition of in Table 2 and the inequality $s \le L$; and the fourth inequality follows from the bound $\delta_{LS} \le 1$. Therefore, since the trust region constraint $\sigma_i \le \frac{v_{i+1}}{s}$ is satisfied, the following holds with probability at least 1 - p:

$$f(\bar{x}) - f(\bar{x}) \le f \quad x - \sigma_i \frac{v_{i+1}}{v_{i+1}} - f(\bar{x}) \le \rho(f(x) - f(\bar{x})).$$

Thus, the first item of the proposition follows.

Contraction case 2: tangent step. Next, we suppose that there exists σ_i satisfying Item 2 of Theorem 4.3. In the interest of analyzing $u_i \in \partial \sigma_i f(x)$, let us show that x, σ_i , and T satisfy the conditions of Proposition 5.6: $x \in B\delta_A/2(\overline{x})$ since $\delta_{LS} \leq \delta_A/2$. Second, by Item 2a of Theorem 4.3, we have

$$\sigma_i \leq a_2 \ y - \overline{x} \leq \ C_5 \ y - \overline{x}.$$

Finally, by Item 2b of Theorem 4.3, we have

$$dist(x, M)/\sigma_i \le 2a_2 \ y - \bar{x} \le C_5 \ y - \bar{x}$$

Therefore, since f(x) - f(x) > f(x) Proposition 5.6 implies that

$$f \quad x - \sigma_i \frac{u_i}{u_i} \quad - f(\overline{x}) \leq f(x) - f(\overline{x}) - \frac{\sigma_i}{8} \operatorname{dist}(0, \, \partial_{\sigma_i} f(x)) \leq \rho(f(x) - f(\overline{x})).$$

Next, we show that u_i and σ_i satisfy the trust region condition $\sigma_i \leq \frac{u_i}{s}$. To show this, we first note that σ_i and x satisfy the conditions of Lemma 4.2: First $x \in B\delta_A/2(x)$ since $\delta_{LS} \leq \delta_A/2$. Second, by Item 2a of Theorem 4.3, we have

$$\sigma_i \leq a_2 \ y - \overline{x} \leq \ C_2 \ y - \overline{x}.$$

Finally, by Item 2 of Theorem 4.3, we have

$$dist(x, M) \le 2a_2\sigma_i \ y - \bar{x} \le 2a_2^2 \ y - \bar{x}^2 \le C_2 \ y - \bar{x},$$

where the third inequality follows from the bounds $y-\overline{x} \le 2\delta_{LS} \le 1/a_2$ and $a_2 \le C_2/2$ (recall that $C_5 \le C_2$). Therefore, by Lemma 4.2 we have $u_i \ge P_{TM}(y)u_i \ge C_1 y-\overline{x}$. Consequently, we have

$$\sigma_i \le a_2 \ y - \overline{x} \le \frac{C_1 \ y - \overline{x}}{s} \le \frac{u_i}{s}$$

where the second inequality follows from the definition of a $_2$ in Table 2 and the inequality $s \le L$. To complete the proof, observe that $V_{i+1} = u_i$: since the sufficient descent condition is met, namely $f(x - \sigma_i u_i / u_i) \le f(x) - \sigma(u_i)$, NDescent terminates at the first iteration. Therefore, we must have

$$f(\widetilde{x}) - f(\overline{x}) \le f \quad x - \sigma_i \frac{V_{i+1}}{V_{i+1}} \quad - f(\overline{x}) \le \rho(f(x) - f(\overline{x})),$$

as desired.

Having proved the desired contraction $f(\bar{x}) - f(\bar{x}) \le \rho(f(x) - f(\bar{x}))$, we now turn to the bound on $\bar{x} - x$.

Stepsize bound. We now no longer assume that $f(x) - f(x) > \max\{2, G, 1, T\}$. We claim that we have

$$\max_{0 \le i \le G - 1} \{ \sigma_i : \sigma_i \le v_{i+1} / s \}$$

$$\le C_6 \max_{1, T} s_{lb}, \ 2 / s_{lb}, \ \overline{2(f(x) - f(x))/ \min\{s_{lb}, \gamma\}} .$$
(6.3)

Note that inequality (6.3) immediately yields the second item of the proposition, since

$$\tilde{x} - x \le \max_{0 \le i \le G - 1} \left\{ \sigma_i : \sigma_i \le v_{i+1} / s \right\}.$$

To prove (6.3), we will apply Lemma 5.8.

To that end, first note that $x \in B\delta_{Grid}(x)$ since $\delta_{LS} \le \delta_{Grid}$. Next, we verify that there exists an index i such that σ_i satisfies a slightly stronger version of the assumptions of Lemma 5.8. Indeed, recall that by the quadratic growth condition (A1), we have the bound

$$\operatorname{dist}(x, M) \le x - \overline{x} \le \overline{2(f(x) - f(\overline{x}))/y}. \tag{6.4}$$

Thus, to satisfy the assumptions of Lemma 5.8, we prove that there exists i such that

$$R_x := D_1^{-1} \quad \overline{2(f(x) - f(\overline{x}))/\gamma} \le \sigma_i \le \delta_{\text{Grid}}. \tag{6.5}$$

Indeed, first notice that $\sigma_0 \le \delta_{Grid}$ since $G \ge \log_2(1/\delta_{Grid})$. Thus, if $\sigma_0 \ge R_x$, the bound (6.5) holds for σ_0 . If instead $\sigma_0 \le R_x$, we have

$$\sigma_0 < R_x \le D_1^{-1} \ \overline{2L \delta_{\mathsf{LS}/V}} \le \ \min\{\delta_{\mathsf{Grid}}/2, 1/4\} \le \min\{\delta_{\mathsf{Grid}}, 1/2\} \le 1/2 = \sigma_{G^{-1}}$$

where the second inequality follows since $x - \overline{x} < \delta_{LS}$ and f is L- Lipschitz continuous on $B_{\overline{\delta}_{LS}}(\overline{x})$; and the third inequality follows since $S_{LS} \leq y$ $D_1^2 \min\{\delta_{Grid}/2, 1/4\}^2/(2L)$. Thus, there exists i such that $\sigma_i \in [\min\{\delta_{Grid}/2, 1/4\}, \min\{\delta_{Grid}, 1/2\}]$. Since $\min\{\delta_{Grid}/2, 1/4\} \geq R_X$, inequality (6.5) follows.

Now let i_* be the minimal such index such that (6.5) is satisfied for $i = i_*$. If $i_* = 0$, the bound $\sigma_{i_*-1} \le R_X$ holds. In particular, $\sigma_{i_*} \le 2R_X$. Therefore, considering the cases $i_* = 0$ and $i_* = 0$ separately, we have

$$R_x \le \sigma_{i_*} \le \max\left\{\sigma_0, 2R_x\right\}. \tag{6.6}$$

Now we bound the step length $x - \tilde{x}$ by considering two cases.

First suppose that $\sigma_{i_*} > u_{i_*}/s$. In this case, (2.1) ensures $\sigma_{i_*} > v_{i_*+1}/s$. Then, since σ_{i} is increasing in i, we have

$$\max_{0 \le i \le G - 1} \{ \sigma_i : \sigma_i \le v_{i+1} / s \} \le \sigma_i,$$

$$\le \max \{ \sigma_0, 2R_x \}$$

$$\le C_6 \max_{2, G} s_{lb}, \quad 2(f(x) - f(\overline{x})) / \min\{s_{lb}, y\} .$$

$$\le C_6 \max_{1, T} s_{lb}, \quad 2(f(x) - f(\overline{x})) / \min\{s_{lb}, y\} .$$

which verifies (6.3). We now consider the alternative case.

Next suppose that $\sigma_{i_*} \le u_{i_*}/s$. We consider two subcases. First suppose that the following bound also holds:

$$u_{i_*} \le \frac{8(f(x) - f(\overline{x}))}{\sigma_{i_*}}.$$
(6.7)

Then, since $\sigma_{i_*} \geq R_x$, we have

$$u_{i_*} \le \overline{32 Y D_1^2(f(x) - f(x))}.$$

Second, suppose that (6.7) does not hold. Let us apply Lemma 5.8 to $\sigma = \sigma_{i_*}$:

$$u_{i_{*}} \leq \max \quad 1 - \frac{\mu^{2}}{64L^{2}} \quad L, 4C_{(a)} \max\{\sigma_{0}, 2R_{x}\} + 4(C_{(a)} + 2\beta) \quad x - \overline{x}$$

$$\leq \max \quad 1 - \frac{\mu^{2}}{64L^{2}} \quad , 8C_{(a)}\sigma_{0}, 8(C_{(a)} + 2\beta) \quad x - \overline{x} + 16C_{(a)}R_{x}$$

$$\leq \max \quad 1 \cdot _{1,T}, 1 \cdot _{2,G}, 8(C_{(a)} + 2\beta + 2C_{(a)}D_{1}^{-1}) \quad \overline{2(f(x) - f(\overline{x}))/y}$$

$$\leq sC_{6} \max\{_{1,T}/s_{lb}, _{2,G}/s_{lb}, \quad \overline{2(f(x) - f(\overline{x}))/y}\},$$

where first inequality follows from Lemma 5.8 and bound (6.6); the second inequality follows from the bound: $\max\{a,b\} + c \le a+b+c \le 2 \max\{a,b+c\}$ for all $a,b,c \ge 0$; the third inequality follows by definition of $_{1,T}$, $_{2,G}$ and R_x , and (6.4); and the last inequality follows since $C_6 \ge \max\{1,8(C_{(a)}+2\beta+2C_{(a)}D_1^{-1})/s_{lb}\}$.

Therefore, as long as $\sigma_{i_*} \leq u_{i_*}/s$, we have

123

$$\begin{aligned} &u_{i_*}/\ s\\ &\leq \max\ C_6 \max\{\ _{1,\,T}/s_{\text{lb}},\ _{2,\,G}/s_{\text{lb}},\ \overline{2(f(x)-f(\overline{x}))/y}\},\quad \overline{32y\,D_1^2(f(x)-f(\overline{x}))/s_{\text{lb}}^2}\\ &\leq C_6 \max\ _{1,\,T}/s_{\text{lb}},\ _{2,\,G}/s_{\text{lb}},\ \overline{2(f(x)-f(\overline{x}))/y}\\ &\leq C_6 \max\ _{1,\,T}/s_{\text{lb}},\ _{2,\,G}/s_{\text{lb}},\ \overline{2(f(x)-f(\overline{x}))/\min\{s_{\text{lb}},\ y\}}\ ,\end{aligned}$$

where second inequality follows from the bound $C_{6} \ge 4V D_{1}/(s_{lb})$; and the third inequality follows from the bound (6.4). To complete the proof of (6.3), recall that by (2.1), for all $j \ge i_*$, we have $V_{j} \le u_{i_*}$. Consequently,

$$\max_{0 \le i \le G - 1} \{ \sigma_i : \sigma_i \le v_{i+1} / s \}$$

$$\le \max_{0 \le i \le G - 1} \sigma_{i_*}, \ u_{i_*} / s$$

$$= u_{i_*} / s$$

$$\le C_6 \max_{1, T} s_{lb}, \ 2, G / s_{lb}, \ \overline{2(f(x) - f(x)) / \min\{s_{lb}, y\}} ,$$

$$F_0 \subset G_0$$

Parameter	Definition
C	$\frac{2048L^2}{\mu^2}$
С	$\max \frac{(C_5^2 L + \beta) L}{C_3}, L, \frac{L}{\min\{1 \cdot a_1\}} + \frac{\beta}{2 \min\{1 \cdot a_1\} a_2^2}, 8C(a)$
q	$\max \ \rho, \ 1 - \frac{3\mu^2}{256L^2}, \frac{1}{2}$
δ_{NTD}	$\min \frac{\delta_{LS}}{4}, \frac{\delta_{LS}^2 \min\{s_{lb}, \gamma\}(1-q^{1/2})^2}{32LC_6^2}, \frac{\delta_{LSs_{lb}}(1-q)}{4LC_6}$
K_0	$\max \ \log_q \ \frac{\delta_{LS}^2 \min\{s_{lb}, \gamma\}(1-q^{1/2})^2}{32CC_6^2} \ , \log_q \ \frac{\delta_{LS}s_{lb}(1-q)}{4CC_6} \ , \log_2 \ \frac{1}{\delta_{Grid}}$

Table 3 Parameters used throughout Sect. 6.3; see also Tables 1 and 2

which verifies (6.3).

6.3 Main Convergence Theorems

We are now ready to prove the main results of this work. The goal of this section is to prove that an event of the following form occurs with high probability.

Definition 6.2 $(E_{k_0,q,C})$ For any $k_0 \ge 0$, $q \in (0,1)$ and $C \ge 0$, let $E_{k_0,q,C}$ denote the event that for all $k \ge k_0$, we have the following two bounds:

$$\begin{split} f(x_k) - & f(\overline{x}) \leq \max\{(f(x_{k_0}) - f(\overline{x}))q^{k-k_0}, Cq^k\}; \\ x_k - & \overline{x}^{-2} \leq \frac{2}{y}\max\{(f(x_{k_0}) - f(\overline{x}))q^{k-k_0}, Cq^k\}. \end{split}$$

We will lower bound the probability of the event L_{0}^{F},q,C in both nonconvex and convex settings for a particular choice of k_0 , q, and C. In the nonconvex setting, our result will lower bound the conditional probability of $E_{k_0,q,C}$, given that iterate x_{k_0} enters a sufficiently small neighborhood of \overline{x} . To prove the result, we will simply iterate the one-step improvement bound of Proposition 6.1. In the convex setting, we will lower bound the unconditional probability of $E_{k_0,q,C}$. To prove this result, we will combine the conditional result with the sublinear convergence guarantee of Theorem 2.4.

Before turning to the proofs, we introduce the main parameters that are common to both the nonconvex and convex settings.

6.3.1 The Nonconvex Setting

The following theorem is our main convergence theorem in the nonconvex setting.

Theorem 6.3 (Main Theorem: Nonconvex Setting) Assume the assumptions outlined at the start of Sect. 6 are satisfied. Recall the notation of Table 3. Fix a failure probability $p \in (0, 1)$ and an index $k_0 \ge \max_{i=1}^{n} K_0$, $C_i \log_i C_i / p_i$. Suppose

 $P(x_{k_0} \in B\delta_{NTD}(\vec{x})) > 0$. Then,

$$P(E_{k_0,q,C} \mid x_{k_0} \in B\delta_{NTD}(\overline{x})) \geq 1 - p$$

Proof We begin with preliminary notation and bounds.

Fix $k_0 \ge \max\{K_0, C \log(C/p)\}$ and for all $k \ge k_0$, define the quantity

$$R_k := \max\{(f(x_{k_0}) - f(\bar{x}))q^{k-k_0}, Cq^k\}.$$

Note that whenever $x_{k_0} \in B\delta_{NTD}(\vec{x})$ we have the bound

$$R_k \le \max\{L \delta_{\mathsf{NTD}} q^{k-k_0}, Cq^k\}, \tag{6.8}$$

since *f* is *L*-Lipschitz continuous on $B \delta(x)$.

Next, we prove that

$$\max\{_{1,T_{k}},_{2,G_{k}}\} \le R_{k+1} \quad \text{for all } k \ge 0.$$
 (6.9)

Indeed, beginning with $_{1,T_{k}}$, we have

$$1_{1,T_{k}} = \max \frac{\left(C_{5}^{2}L + \beta\right)L}{C_{3}} \quad 1 - \frac{3\mu^{2}}{256L^{2}} \quad , \quad 1 - \frac{\mu^{2}}{64L^{2}} \quad L$$

$$\leq C \max \left\{ 1 - \frac{3\mu^{2}}{256L^{2}} \quad , \quad 1 - \frac{\mu^{2}}{64L^{2}} \right\}$$

$$\leq Cq^{k+1} \leq R_{k+1},$$

where the first and second inequalities follow from the definitions of C and q together with the lower bound $T_k \ge k + 1$. Turning to $_{2,G_k}$, we have

$$\begin{array}{ll} _{2,G_{k}}=\max & \frac{L}{\min \{1,\,a_{1}\}}+\frac{\beta}{2\,\min \{1,\,a_{1}\}a_{2}^{2}},\,8C(a),\,L & 2^{-\,G_{k}}\leq C2^{-\,G_{k}}\\ \leq Cq^{\,k+\,1}\leq \,R_{k+\,1}, \end{array}$$

where the first and second inequalities follow from the definition of C and q together with the lower bound $G_k \ge k + 1$. Thus (6.9) holds.

Finally, we analyze the quantity

$$D_{k_0}, \delta_{\text{NTD}} := \int_{k=k_0}^{\infty} C_6 \max \frac{2R_k/V}{r}, R_{k+1}/s_{\text{lb}} \quad \text{where } V := \min\{s_{\text{lb}}, y\}.$$

We claim in particular that

$$D_{k_0}, \delta_{\text{NTD}} + \delta_{\text{NTD}} \le \delta_{\text{LS}}/2. \tag{6.10}$$

Since $\delta_{NTD} \leq \delta_{LS}/4$, it suffices to prove $D_{k_0}, \delta_{NTD} \leq \delta_{LS}/4$. To that end, we have

$$\begin{split} D_{k_0}, & \delta_{\mathsf{NTD}} = \int_{k=k_0}^{\infty} C_6 \max \frac{\overline{2R_k/V}}{2 \operatorname{max} \{L \delta_{\mathsf{NTD}} q^{k-k_0}, Cq^k\}/V}, \max\{L \delta_{\mathsf{NTD}} q^{k-k_0}/V, Cq^k\}/s_{\mathsf{lb}} \\ & \leq C_6 \max \frac{2 \operatorname{max} \{L \delta_{\mathsf{NTD}} q^{k-k_0}, Cq^k\}/V}{\overline{V} (1-q^{1/2})}, \frac{\overline{2Cq^{k_0}}}{\overline{V} (1-q^{1/2})}, \frac{L \delta_{\mathsf{NTD}}}{s_{\mathsf{lb}} (1-q)}, \frac{Cq^{k_0}}{s_{\mathsf{lb}} (1-q)} \\ & \leq \frac{\delta_{\mathsf{LS}}}{4}, \end{split}$$

where the first inequality follows from the bounds (6.8) and the bound $R_{k+1} \le R_k$; the second inequality follows by summing the infinite series; and the third inequality follows from the definitions of K_0 and δ_{NTD} together with the bound $k_0 \ge K_0$. This proves (6.10).

We now turn to the proof. Consider the following sequence defined for all $k \ge k_0$:

$$b_k := \delta_{NTD} + \sum_{j=k_0}^{k-1} C_6 \max \frac{2R_k/y}{R_{k+1}/s_{lb}}$$
.

Note that (6.10) ensures that $b_k \le \delta_{LS}/2$ for all $k \ge k_0$. Now, define the event

$$F_{k_0} := \{ x_{k_0} \in B_{\delta_{NTD}}(\bar{x}) \}.$$

In addition, define the following decreasing sequence of events

$$A_k := \sum_{j=k_0}^{2^k} f(x_j) - f(\bar{x}) \le R_j \text{ and } x_j - \bar{x} \le b_j$$
.

We claim that

$$P(A_{k+1} \mid A_k \cap F_{k_0}) \ge 1 - \exp(-T_k/C)$$
 for all $k \ge k_0$. (6.11)

Indeed, Proposition 6.1 implies that conditioned on $A_k \cap F_{k_0}$, the following four inequalities are satisfied with probability at least $1 - \exp(-T_k/C)$:

1.
$$f(x_k) - f(\bar{x}) \le R_k$$

2. $x_k - \bar{x} \le b_k$;

2.
$$x_k - x \leq b_k$$
;

3.
$$x_{k+1} - x_k \le C_6 \max_{1, T_k} / s_{lb}, \frac{1}{2 \cdot G_k} / s_{lb}, \frac{1}{2 \cdot f(x_k) - f(\overline{x})} / Y$$
;
4. $f(x_{k+1}) - f(\overline{x}) \le \max\{\rho(f(x_k) - f(\overline{x})), \frac{1}{1 \cdot T_k}, \frac{2}{2 \cdot G_k}\}.$

(Note that in applying the Proposition 6.1, we use the scalar $s = \max\{g_k, c_0 g_0\}$ and the inclusion $s \in [s_{lb}, L]$, which was proved (6.1).) Thus, the bound (6.11) will follow by induction if we can prove that whenever the above four inequalities hold, we have $x_{k+1} - \overline{x} \le b_{k+1}$ and $f(x_{k+1}) - f(\overline{x}) \le R_{k+1}$.

To that end, we first prove $x_{k+1} - \overline{x} \le b_{k+1}$. Indeed,

$$x_{k+1} - \overline{x} \le x_{k+1} - x_k + x_k - \overline{x}$$

$$\le C_6 \max_{1, T_k} / s_{lb}, \ 2 \cdot G_k / s_{lb}, \ \overline{2(f(x_k) - f(\overline{x}))/y} + b_k$$

$$\le C_6 \max_{1, T_k} \overline{2R_k/y}, R_{k+1} / s_{lb} + b_k = b_{k+1}, \tag{6.12}$$

where the second inequality follows from Proposition 6.1; and the third inequality follows from the bound (6.9). Next, we prove the bound on $f(x_{k+1}) - f(x) \le R_{k+1}$. Indeed,

$$f(x_{k+1}) - f(\overline{x}) \leq \max\{\rho(f(x_k) - f(\overline{x})), x_{1,T_k}, x_{2,G_k}\}$$

$$\leq \max\{\rho \max\{(f(x_{k_0}) - f(\overline{x}))q^{k-k_0}, Cq^k\}, x_{1,T_k}, x_{2,G_k}\}$$

$$\leq R_{k+1},$$

where the final inequality follows from (6.9) and the bound $\rho \leq q$. Consequently, the bound (6.11) holds. Moreover, due to the bound $T_k \geq k + 1$, we have

$$P(A_{k+1} \mid A_k \cap F_{k_0}) \ge 1 - \exp(-T_k/C) \ge 1 - \exp(-(k+1)/C).$$
 (6.13)

Now we relate A_k to $E_{k_0,q,C}$. To that end, by the conditional law of total probability, for all $k \ge k_0$, we have

$$P(A_{k+1} \mid F_{k_0}) \geq P(A_{k+1} \mid A_k \cap F_{k_0}) P(A_k \mid F_{k_0}) \geq P(A_k \mid F_{k_0}) - \exp(-(k+1)/C).$$

Therefore, for all $k \ge k_0$, we have

$$P(A_k \mid F_{k_0}) \ge P(A_{k_0} \mid F_{k_0}) - \sum_{j=k_0+1}^{\infty} \exp(-j/C) = 1 - \frac{\exp(-\frac{k_0+1}{C})}{1 - \exp(-\frac{1}{C})} \le 1 - p,$$

where the equality follows since $P(A_{k_0} | F_{k_0}) = 1$; and the final inequality follows by definition of $k_0 \ge C \log(C/p)$. Now recall that $\sup_{k \ge k_0} b_k \le \delta_{LS}/2$. Therefore, defining the event

we have

$$P(E_{k_0,q,C} \mid F_{k_0}) \ge \lim_{k \to \infty} P(A_k \mid F_{k_0}) \ge 1 - p$$

Next, recall that since $\delta_{LS} \leq \delta_A$, the quadratic growth bound (A1)

$$x_k - \overline{x}^{-2} \le \frac{2}{V} (f(x_k) - f(\overline{x})) \le \frac{2}{V} R_k$$

holds for every $k \ge k_0$ within the event $E_{k_0,q,C}$. Thus, $E_{k_0,q,C} \supseteq E_{k_0,q,C}$. Therefore, we have

$$P(E_{k_0,q,C} \mid F_{k_0}) \ge P(E_{k_0,q,C} \mid F_{k_0}) \ge 1 - p$$

as desired.

6.3.2 The Convex Setting

Now we turn to the convex setting. Our goal is to prove a lower bound on $R(E_{k_0,q,C})$ for q and C chosen as in Table 3 and all sufficiently large k. Before stating the result, we recall a simple fact about convex functions satisfying Assumption A. A similar result appears in [8, Section 2.4], but for completeness we provide a proof in Appendix F.

Lemma 6.4 In addition to the assumption set out at the start of the section, suppose that function f is convex. Then for all a > 0, we have

$$\{x \in \mathsf{R}^d : f(x) - f(\bar{x}) \leq a\} \subseteq \overline{B}_{r_a}(\bar{x}) \qquad where \ r_a := \max \ \frac{2a}{y \ \delta_{\!\!A}}, \ \frac{2a}{y} \ .$$

In particular, f has bounded sublevel sets.

We now turn to our main theorem.

Theorem 6.5 (Main Theorem: Convex setting) Assume the assumptions of Sect. 6.1 are satisfied. Recall the notation of Table 3. In addition, suppose that function f is convex. Consider the bounded set

$$S:=\{x+u: f(x)\leq f(x_0) \ and \ u\in \overline{B}(x)\}.$$

Let L be a Lipschitz constant of f on S. Define the constants

$$a := \min \quad \frac{\gamma \, \delta_{\!\!A} \delta_{\!\! \text{NTD}}}{4}, \frac{\gamma \, \delta_{\!\!\! \text{NTD}}^{\!\!\! \text{NTD}}}{8} \quad ; \qquad and \qquad b := \inf_{\alpha \in (0,1)} \frac{64L \quad \frac{2}{\alpha}}{\frac{a}{\text{diam}(S)}} \frac{\frac{2\alpha}{(1-\alpha)}}{\frac{a}{\text{diam}(S)}}.$$

Finally, define

$$K_1 := \quad \frac{4 \mathrm{diam}^2(S)}{a^2} \min \ \ 16^2(f(x_0) - \inf f)^2, \\ \frac{b}{4}, 2048L^2 \log \ \ \frac{2}{p} \ \ , 128(L)^2 \ + \frac{(4L)^2}{a^2} \ \ .$$

Then, for every failure probability $p \in (0, 1)$, we have

$$P(E_{k_0,q,C}) \geq 1 - p \quad \textit{ for all } k_0 \geq \max \ K_0, C \ \log \ \frac{2C}{p} \ , 2K_1 - 1 \ .$$

Proof Theorem 6.3 shows that

$$P(E_{k_0,q,C} \mid x_{k_0} \in B\delta_{NTD}(\overline{x})) \ge 1 - p/2$$
 for all $k_0 \ge \max K_0$, $C \log \frac{2C}{p}$ (6.14)

We claim that

$$P(x_{k_0} \in B_{\delta_{NTD}}(\overline{x})) \ge 1 - p/2 \quad \text{for all } k_0 \ge 2K_1 - 1.$$
 (6.15)

Note that this yields the proof, since in that case

$$P(E_{k_0,q,C}) \ge P(E_{k_0,q,C} \mid x_{k_0} \in U) P(x_{k_0} \in B\delta_{NTD}(\overline{x})) \ge 1 + p^2/4 - p \ge 1 - p$$

for all $k_0 \ge \max K_0$, $C \log \frac{2C}{p}$, $2K_1 - 1$.

Observe that (6.15) will follow if

$$P(f(x_{k_0}) - f(\bar{x}) \le a) \ge 1 - p/2$$
 for all $k_0 \ge 2K_1 - 1$. (6.16)

Indeed, by Lemma 6.4, we have.

$$\{x \in \mathbb{R}^d : f(x) - f(\overline{x}) \le a\} \subseteq \overline{B} \delta_{NTD}/2(\overline{x}) \subseteq B \delta_{NTD}(\overline{x}).$$

To prove (6.16), we apply Theorem 2.4. To that end, note $t[tat \in \mathbb{R}^d : f(x) \le f(x_0)]$ and the widened sublevel set S are indeed bounded, due to Lemma 6.4. Therefore D and the Lipschitz constant L of f on S are finite. Now observe $G := \min_{K_1 \le k \le 2K_1 - 1} \{G_k\} \ge K_1$ since $G_k \ge k + 1$ for all k. Thus, there exists $i \le G$ such that

$$(1/2)K_1^{-1/2} \le \sigma_i \le K_1^{-1/2}$$

Therefore, applying Theorem 2.4 (in particular (2.2)) with this σ_i , we have

$$f(x_{2K_1-1}) - f(x)$$
23 $\lim_{n \to \infty} \frac{f(x)}{n}$

$$\leq D \max \left\{ \frac{16(f(x_{K_1}) - \inf f)}{K_1^{1/2}}, \frac{16L}{K_1^{1/2}}, \frac{2\log(2K_1^{2/}p)}{K_1^{1/2}}, \frac{\sqrt{\frac{128}{128}L}}{K_1^{1/2}} \right\} + \frac{2L}{K_1^{1/2}}$$
(6.17)

with probability at least 1 - p/2. Thus, to complete the proof, we show that the left-hand side of (6.17) is smaller than a. Indeed, it is straightforward to check that

$$\max \ \frac{2L}{K_1^{1/2}}, \frac{16D(f(x_{K_1}) - \inf f)}{K_1^{1/2}}, \frac{\sqrt[4]{128}DL}{K_1^{1/2}} \le \frac{a}{2}.$$

Thus, the proof will follow if

$$\frac{16DL}{K_1^{1/2}} \frac{2\log(2K_1^{2/}p)}{\leq \frac{a}{2}}.$$
 (6.18)

We perform this calculation in Appendix G. Thus, the proof is complete.

6.3.3 Local Oracle Complexity

Thus, we have established a local nearly linear convergence rate for NTDescent. To understand the overall complexity of the method, we must derive an upper bound on the contraction factor q. The following lemma, which is proved in Appendix H, provides one that depends on a worst-case condition number of f.

Lemma 6.6 Suppose that without loss of generality that $\delta_A \leq 1$. Define the condition number

$$\kappa = \frac{\max\{L, \beta, C_{(a)}\}}{\min\{y, \mu\}}.$$

Then there exists a universal constant $\eta > 0$ independent of f such that

$$q \leq 1 - \frac{\eta}{\kappa^8 (1 + C_M)^2}.$$

where q is defined as in Table 3.

With this upper bound on q, it is straightforward to derive a local complexity estimate for NTDescent: the method locally produces a point \hat{x} satisfying $f(\hat{x}) - f(\bar{x}) \le \varepsilon$ with at most

$$O = \kappa^8 (1 + C_M)^2 \log(1/\epsilon)^{-3}$$
,

first-order oracle evaluations. This bound may be pessimistic since we did not attempt to optimize the constants C_i or a_i . We leave the improvement of this complexity as an intriguing open question.

Before moving to a brief numerical illustration, we explain how Theorem 1.1 from the introduction follows from the above results.

Remark 2 (Establishing Theorem 1.1) Theorem 1.1 from the introduction immediately follows from Theorems 6.3 and 6.5. Indeed, first the event $E_{q_0,q,C}$ from Theorems 6.3 and 6.5 contains the corresponding event $E_{k_0,q,C}$ from Theorem 1.1 for particular q and C, which depend solely on f. Second, from the statement of theorems, we see that the neighborhood of local nearly linear convergence, $E_{NTD}(\vec{x})$, depends solely on f.

7 Numerical Illustration

In this section, we briefly illustrate the numerical performance oNTDescent on two nonsmooth objective functions, borrowed from [1, 12, 37, 39]. In both experiments, we compare NTDescent to the subgradient method with the popular Polyak stepsize (PolyakSGM) [47], which iterates

$$x_{k+1} = x_k - \frac{f(x_k) - \inf f}{W_k^2} W_k$$
 for some $W_k \in \partial f(x_k)$.

In the first example, inf f is known, in the second, we estimate inf f from multiple runs of NTDescent. We choose to compare against the subgradient method because it is a simple first-order method with strong convergence guarantees in convex [47] and nonconvex settings [20]. Importantly, PolyakSGM accesses the objective solely through function and subgradient evaluations. Thus, we compare the accuracy achieved by PolyakSGM and NTDescent after a fixed number of oracle calls, i.e., evaluations of ∂f .

Let us comment on the implementation of NTDescent. First, in all experiments, unless otherwise noted, we do not tune parameters of TDescent. Instead, we simply choose scaling constant $c_0 = 10^{-6}$ and loop size parameters

$$T_k = k + 1$$
 and $G_k = \min\{k + 1, \log_2(10^{-16})\}$ for all $k \ge 0$.

Second, we attempt to save first-order oracle calls by breaking the loop on Lines 2 through 6 of Algorithm 3 whenever we find tha $\sigma_i > v_{i+1}/s$. Since σ_i is increasing in i and v_{i+1} is nonincreasing in i, this does not affect the iterates of NTDescent; see Lemma 2.1. Finally, in all problems, we initialize NTDescent and PolyakSGM at a random vector az where z is sampled from the uniform distribution on the unit sphere. For all problems, we use a=1 unless otherwise noted. Note that in the problems of Sect. 7.1 and 7.2, the solution is known, while in the problem of Sect. 7.3, the solution is unknown.

The purpose of this section is not to argue that NTDescent is a substitute for standard subgradient methods in most problems. Instead, we only wish to point out

some scenarios where standard first-order methods are known to perform poorly, yet NTDescent asymptotically accelerates. We are also not arguing that NTDescent has fast global rates: indeed, we previously mentioned that the NTDescent's global rate is $O(^{-6})$ which is much worse than PolyakSGM's $O(^{-2})$ rate for general convex problems. In practice, one could devise schemes that couple NTDescent with PolyakSGM, eventually switching toNTDescent when it begins to outperform PolyakSGM. While we leave a more thorough numerical study to future work, the reader may download and run our PyTorch [45] implementation of NTDescent at the following url: https://github.com/COR-OPT/ntd.py

We now turn to the examples.

7.1 A Max-of-Smooth Function

In this example, *f* takes the following form

$$f(x) = \max_{i=1,\dots,m} g_i x + \frac{1}{2} x^T H_i x , \qquad (7.1)$$

where we generate a random vector $\lambda \in \mathbb{R}^m$ in $\{\lambda > 0 : \prod_{i=1}^m \lambda_i = 0\}$, a random positive semi-definite matrix H, and a random vector gsatisfying that $\prod_{i=1}^m \lambda_i g_i = 0$. In this case, one can show that with probability \mathcal{Y}_i , satisfies Assumption A at its unique minimizer 0.

In Fig. 6 we plot the performance of NTDescent and PolyakSGM for multiple pairs of (d, m), varying initialization scale, a slight modification of NTDescent that allows longer steps, and varying scales c_0 . We begin with Fig. 6a, b. Figure 6a shows that the performance of NTDescent depends on m. On the other hand, Fig. 6b shows NTDescent performance is independent of d, as expected. Both plots show that NTDescent asymptotically outperforms PolyakSGM· Turning to initialization, Fig. 6c shows the result of initializing NTDescent at a random vector az, where z is uniformly drawn from the sphere and a is a scale parameter satisfying $a\{1, 10, 100\}$. Clearly, NTDescent is affected by the initialization scale but surpasse PolyakSGM after 30,000 oracle calls. While we expect NTDescent to converge slowly when far from minimizers, we introduce a simple strategy to mitigate this behavior.

Adaptive grid strategy. Briefly, suppose we run linesearch the full G steps without exiting (via the violation of the trust region constraint). Then we simply continue the linesearch loop trying $\sigma_{-1} = 10\sigma_0$, $\sigma_{-2} = 10\sigma_{-1}$, and so on, until we violate the trust region constraint or σ_{-i} exceeds a predefined threshold.

Figure 6d shows the result of this strategy with a predefined threshold ∞ , showing that it compensates for poor initialization quality. Finally in Fig. 6e, f, we show the effect of changing the c_0 input to NTDescent. It appears NTDescent is relatively insensitive to c_0 and smaller choices generally result in better performance. This motivates our default choice $c_0 = 10^{-6}$ in the remainder of the experiments.

Before turning to our second experiment, we briefly mention two alternative methods—Prox-linear [13, 27, 28, 51, 54] and Survey Descent [30]—which could

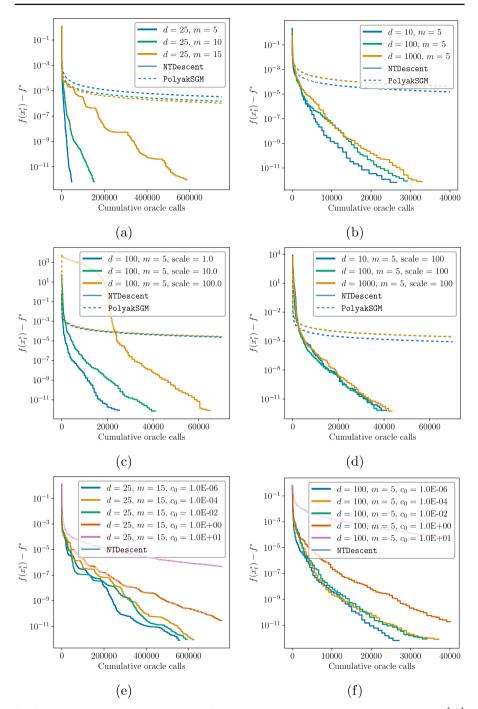


Fig. 6 Comparison of NTDescent with PolyakSGM on (7.1). For both algorithms, the value $f(x_t^*)$ denotes the best function seen after t oracle evaluations. See text for description

be applied to this problem. In order to explain these algorithms, let us write $f = \max_{i=1,...,m} \{f_i\}$, where the f_i are the quadratic function from (7.1).

Prox-linear method. Given a point $x \in \mathbb{R}^d$, the Prox-linear update x_+ solves

$$x_{+} = \underset{y \in \mathbb{R}^{d}}{\operatorname{argmin}} \max_{i = 1, \dots, m} \{ f_{i}(x) + \nabla f_{i}(x), y - x \} + \frac{\rho}{2} y - x^{2}.$$

One may show that x_+ geometrically improves on x; see [25]. However, in contrast to NTDescent, the prox-linear method requires that the components f_i are known. This is stronger than the first-order oracle model considered in this work. Thus, we do not compare NTDescent with prox-linear.

Survey Descent The Survey Descent method is a multi-point generalization of gradient descent, designed for max-of-smooth functions. Rather than maintaining a single iterate sequence, the Survey Descent maintains a *survey S* of points, meaning a collection of points $\{s_i\}_{i=1}^m$ at which f is differentiable. A single iteration of the Survey Descent method then aims to produce a new survey $S^+ = \{s_i^+\}_{i=1}^m$ satisfying

$$s_i^+ := \underset{x \in \mathbb{R}^d}{\operatorname{argmin}} \frac{3}{3} x - s_i - \frac{1}{L} \nabla f(s_i) \frac{3}{3}^2$$

$$\text{subject to: } f(s_j) + \nabla f(s_j), x - s_j$$

$$+ \frac{L}{2} |x - s_j|^2 \le f(s_i) + \nabla f(s_i), x - s_i \quad \forall j = i.$$

Here, L is an upper bound on the Lipschitz constant of ∇f_i for all $i = 1, \dots, m$. In [30], Han and Lewis study linear convergence of Survey Descent on max-of-smooth functions under the conditions of Corollary 3.4. Given a survey S, they show that the updated survey S^+ geometrically improves on S (in an appropriate sense) whenever the following conditions are satisfied: (i) all elements of the survey S are near \overline{x} ; (i) the survey S is valid, meaning there exists a permutation a on [m] such that

$$f_{a(i)}(s_i) = f(s_i)$$
 and $\partial f(s_i) = \{ \nabla f_{a(i)}(s_i) \}$ for all $i = 1, \dots, m$.

To estimate the number of components m and find a valid initial survey S sufficiently close to x, Han and Lewis suggest an empirical procedure based on running a nonsmooth variant of BFGS [39] for several iterations. After running BFGS, they suggest to (i) compute an estimate \hat{m} of m from a singular value decomposition of the computed gradients, and (ii) build the survey from \hat{m} past iterates in such a way that the computed gradients form an affine independent set. From the numerical illustration in [30], Survey Descent performs well on several small problems. However, since the initialization procedure and implementation of Survey Descent are somewhat sophisticated, we leave a detailed comparison between NTDescent and Survey Descent and to future work.

7.2 A Matrix Sensing Problem

In this example, f takes the following form

$$f(X) = \frac{1}{n} A(XX^T) - A(M)_{1}$$

where $M \in \mathbb{R}^{N \times N}$ is an unknown positive semidefinite matrix of rank rthat we wish to recover from known linear measurements A(M); the linear operator $A: \mathbb{R}^{N \times N} \to \mathbb{R}^n$ takes the form $Y \to (a_i^T Y a_i - b_i^T Y b_i)_{i=1}^n$, for $n \in N$, where $a_i, b_i \in \mathbb{R}^d$ are random vectors sampled from a standard multivariate normal distribution; and the decision variable is a tall and skinny matrix $X \in \mathbb{R}^{N \times r}$, where in general we allow r = r. This optimization problem appears in various signal processing applications and is known as quadratic sensing [15]. Note that this objective does not satisfy Assumption A, since the solution set is not isolated.

We consider two settings in this section: the exact setting r = r and the overparameterized setting r > r. In the exact setting [14] showed that if n = Nr (), the objective f is sharp, meaning f(x) = d(st(x, argmin f)) and that PolyakSGM converges linearly whenever the initial iterate is sufficiently close to the set of minimizers. In the overparameterized setting, we are not aware of similar guarantees. Note that in practice, r is unknown, so the overparameterized setting is likely to be encountered.

In Fig. 7 we plot the performance of NTDescent and PolyakSGM in two experiments. In Fig. 7a we use base dimensions N=100, optimal rank r=5, and varying overparameterization $r\in\{r,r+2,r+5\}$. In Fig. 7b we use base dimensions N=100, varying optimal rank $r\in\{5,10,15\}$, and fixed overparameterization r=r+5. In both experiments, we fix n=4Nr. Note that the dimension of the decision variable X is varying across each run since d=Nr and r are varying. As can be seen from the plot, PolyakSGM outperforms NTDescent in the exact setting. This is to be expected, since f is a sharp function on which PolyakSGM is known to perform well. On the other hand, when r>r, we find that both methods slow down. However, NTDescent continues to converge nearly linearly, while PolyakSGM converges sublinearly.

7.3 An Eigenvalue Product Function

In this example, we aim to optimize a function \tilde{f} that takes the following form

$$\tilde{f}(X) = \log E_K(A " X),$$

where A is a fixed positive semi-definite data matrix, $E_K(Y)$ denotes the product of K largest eigenvalues of a symmetric matrix $Y \in S^N$, and " denotes the Hadamard (entrywise) matrix product, subject to the constraint that X is positive semi-definite and its diagonal entries are 1. This example is a nonconvex relaxation of an entropy minimization problem arising in an environmental application [1, 12]. In our experiments, we choose A as in [1]: A is the leading $N \times N$ submatrix of a 63 \times 63 covariance

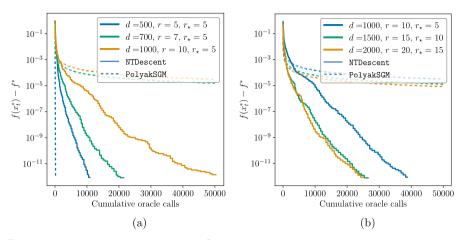


Fig. 7 Comparison of NTDescent with PolyakSGM on (7.1). In both plots, the base dimension is N = 100. Left: fixed optimal rank r = 5 and varying overparameterization $r \in \{5, 7, 10\}$; Right: varying $r \in \{5, 10, 15\}$, fixed overparameterization r = r + 5. For both algorithms, the value $f(x_t^*)$ denotes the best function seen after t oracle evaluations

matrix, scaled so that the largest entry is 1. As suggested by [12], we reformulate this problem as an unconstrained optimization problem using a Burer-Monteiro type factorization

$$\min_{V \in \mathbb{R}^{N \times N}} f(V) = \tilde{f}(c(V)c(V)), \tag{7.2}$$

where $c: \mathbb{R}^{N \times N} \to \mathbb{S}^{N}$ satisfies $c(V) = \operatorname{Diag}([\operatorname{diag}(VV)]^{-1/2})V$ for all $V \in \mathbb{R}^{N \times N}$. Here, the mapping diag (\cdot) takes a matrix an $N \times N$ matrix A to the N dimensional vector with i th entry A_{i} . On the other hand, the mapping $\operatorname{Diag}(\cdot)$ takes an N dimensional vector V to the $N \times N$ diagonal matrix with i th diagonal entry V_i . A formula for the subgradient of f may be found [12]. We do not attempt to verify that f satisfies the full Assumption A. Instead, we point out that under a "transversality condition," function f admits an active manifold at local minimizers [39].

Turning to the experiment, we consider the case where N=14 and K=7. In this example, the optimal function value inf f is not known. Thus, we run NTDescent from four random initial starting points. We terminate each run on NTDescent when a certain "optimality gap" R_k satisfies $R_k \le 10^{-12}$. We denote the minimal function value achieved across all four runs by f. Let us now define and motivate the optimality gap. For iteration k in Algorithm 4, define

$$R_k := \min \max\{\sigma_i^{(k)}, v_{i+1}^{(k)}\}: \sigma_i^{(k)} \le v_{i+1}^{(k)}$$

where $\sigma_i^{(k)}$ and $V_{i+1}^{(k)}$ are computed in Lines 2 through 6 of Algorithm 3 at iteration k. Provided that x_k is sufficiently close to a point \overline{x} at which function f satisfies Assumption A, it is possible to show that R_k satisfies $f(x_k) - f(\overline{x}) - R_k$. This is illustrated in Fig. 8a: there, the optimality gap closely tracks the estimated function

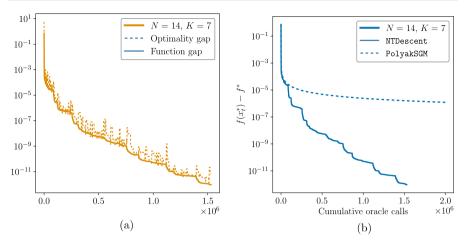


Fig. 8 Numerical performance on (7.2). Left: the close relationship between "optimality gap" and function gap; Right: comparison of PolyakSGM and NTDescent from three initial starting points. For both algorithms, the value $f(x_t^*)$ denotes the best function seen after t oracle evaluations. See text for detail

gap, when approximating by inf f by f *. In Fig. 8b, we compare the performance of NTDescent on the three runs which did not achieve function value f * before termination. In all three cases, we see similar performance. Next, for each run of NTDescent, we also runPolyakSGM from the same initial starting point, estimating inf f by f *. We see that NTDescent outperforms PolyakSGM.

A Proof of Lemma 2.2

Let g denote the minimal norm element of $\partial_{\sigma} f(x)$. Write g as a convex combination of subgradients: $g = \prod_{i=1}^{n} \lambda_i g_i$ where $\prod_{i=1}^{n} \lambda_i = 1$ and $g_i \in \partial f(x_i)$ for some $x_i \in B_{\sigma}(x)$ and $n \ge 0$. Then

$$f(x) \leq f \qquad \lambda_{i}x_{i} + \sum_{i=1}^{n} \lambda_{i}(x - x_{i})$$

$$i = 1 \qquad i = 1$$

$$\leq \sum_{i=1}^{n} \lambda_{i} f(x_{i}) + L\sigma$$

$$i = 1$$

$$\leq f(y) + \sum_{i=1}^{n} \lambda_{i}g_{i}, x_{i} - y + L\sigma$$

$$\leq f(y) + g, x - y + \sum_{i=1}^{n} \lambda_{i} g_{i}, x_{i} - x + L\sigma$$

$$\leq f(y) + \operatorname{dist}(0, \partial_{\sigma} f(x)) x - y + 2L\sigma,$$

as desired.

B Proof of Proposition 3.5

We begin with preliminary notation and bounds. First, since M is C^4 smooth, the projection P_M is C^3 smooth near \bar{x} . Second, since f is C^3 smooth along M near \bar{x} , the composition $f_M := f \circ P_M$ is also C^3 smooth near \bar{x} . Third, the constant μ is positive due to the active manifold assumption. Fourth, choose $\delta > 0$ small enough that the following hold:

- 1. ∇PM is CM -Lipschitz on $B\delta(x)$;
- 2. ∇ *fM* is β -Lipschitz on $B\delta(\vec{x})$;
- 3. $\nabla^2 f M$ is ρ -Lipschitz on $B\delta(\vec{x})$ in the operator norm, where $\rho := 2 \text{lip}_{\nabla^2 f M}^{\text{op}}(\vec{x})$;
- 4. f is L-Lipschitz on $B \delta(x)$;
- 5. the quadratic growth bound (A1) holds:

$$f(x) - f(\bar{x}) \ge \frac{\gamma}{2} |x - \bar{x}|^2$$
 for all $x \in \overline{B}\delta(\bar{x})$;

6. the strong (a) bound (A3) holds:

$$P_{T_M(y)}(v - \nabla_M f(y)) \le C_{(a)} x - y \tag{B.1}$$

for all $x \in \overline{B}\delta(x)$, $v \in \partial f(x)$, and $y \in M \cap \overline{B}\delta(x)$.

7. the $(b \le)$ regularity bound (A4) holds:

$$f(x) \ge f(x) + V, x - x - \frac{\mu}{2} x - \hat{x}$$
 (B.2)

for all $x \in B\delta(\overline{x})$, $v \in \partial f(x)$, and $x \in B\delta(\overline{x}) \cap M$.

8. the sharpness condition holds:

$$\operatorname{dist}(0, \partial f(x)) > 2\mu$$
 for all $x \in B\delta(\overline{x}) \setminus M$.

Given these bounds, let us define

$$\delta_{\rm A} := \frac{1}{2} \min \ \delta, \frac{9 \gamma}{16 \rho}, \frac{\mu}{2 (C_{(a)} + 2 \beta + 2 C_M L)}$$

For this choice of δ_A , Item 1 holds automatically. We now prove the remaining items.

B.1 Item 2: Smoothness of P M

Fix $x \in B_2\delta_A(\vec{x})$ and $x \in B\delta_A(x)$. Observe that $P_M(x) \in B_2\delta_A(\vec{x})$ and we have the inclusion $x - P_M(x) \in N_M(P_M(x))$. Consequently, we have

- 1. $P_{T_M (P_M (x))}(x) = P_{T_M (P_M (x))}(P_M (x));$ 2. $P_M (x) = P_M (P_M (x));$
- 3. $\nabla P_M(P_M(x)) = P_{T_M(P_M(x))}$.

Therefore, we have

$$P_{M}(x) - P_{M}(x) - P_{T_{M}(P_{M}(x))}(x - x)$$

$$= P_{M}(x) - P_{M}(P_{M}(x)) - \nabla P_{M}(P_{M}(x))(x - P_{M}(x))$$

$$\leq \frac{C_{M}}{2} x - P_{M}(x)^{2}$$

$$\leq C_{M}(x - x^{2} + \operatorname{dist}^{2}(x, M)),$$

where the first inequality follows from Lipschitz continuity of ∇PM on $B_2\delta_A(\overline{x}) \subseteq B\delta(\overline{x})$.

B.2 Item 3: Bounds on $\nabla M f$

Recall that $P_M(x) \in B_2\delta_A(x)$ whenever $x \in B\delta_A(x)$. Thus, below we prove that

$$\frac{\gamma}{2} y - x \le \nabla$$
 $f_M(y) \le \beta$ $y - x$ for all $y \in B_2 \delta_A(x) \cap M$.

This is equivalent to the claimed bound since $\nabla fM(y) = \nabla M f(y)$ for all $y \in B_2 \delta_A(\overline{x}) \cap M$.

Let us first prove the claimed upper bound. Due to the inequality,

$$f_M(x) - f_M(\overline{x}) \ge \frac{y}{2} P_M(x) - \overline{x}^2$$
 for all $x \in B\delta(\overline{x})$,

it follows that \overline{x} is a local minimizer of fM. Consequently, $\nabla fM(\overline{x}) = 0$. Thus, since β is a local Lipschitz constant of ∇fM on $B\delta(\overline{x})$, we have

$$\nabla$$
 $f_M(y) \le \beta$ $y - \bar{x}$ for all $y \in B_{\delta}(\bar{x}) \cap M$.

Since $2\delta_A \le \delta$, this proves the claimed upper bound.

Next, we prove the claimed lower bound. It suffices to establish the following convexity inequality:

$$f_M(y) + \nabla f_M(y), \overline{x} - y \le f_M(\overline{x})$$
 for all $y \in B_2 \delta_A(\overline{x}) \cap M$. (B.3)

Indeed, if this inequality holds, we have

$$\nabla f_M(y), y - \overline{x} \ge f_M(y) - f_M(\overline{x}) \ge \frac{\gamma}{2} y - \overline{x}^2 \quad \text{for all } y \in B_2 \delta_A(\overline{x}) \cap M,$$

and the desired result follows from Cauchy-Schwarz.

To that end, observe that since $\nabla f_M(\vec{x}) = 0$ and $\nabla^2 f_M$ is ρ -Lipschitz in $B_2 \delta_A(\vec{x})$, we have

$$f_{M}(y) \leq f_{M}(\bar{x}) + \frac{1}{2} \sqrt[4]{2} f_{M}(\bar{x})(y - \bar{x}), y - \bar{x} + \frac{\rho}{6} y - \bar{x}^{3} \quad \text{for all } y \in B_{2}\delta_{A}(\bar{x}).$$

$$123 \qquad \text{for all } y \in B_{2}\delta_{A}(\bar{x}).$$

Consequently, we have the lower bound on the quadratic form: for $al \not\equiv y B_2 \delta_A(\vec{x}) \cap M$, we have

$$\frac{1}{2} \stackrel{4}{\nabla^{2}} f_{M} (\bar{x})(y - \bar{x}), (y - \bar{x}) \stackrel{5}{\geq} f_{M} (y) - f_{M} (\bar{x}) - \frac{\rho}{6} y - \bar{x}^{3}$$

$$\geq \frac{\gamma}{2} y - \bar{x}^{2} - \frac{\rho}{6} y - \bar{x}^{3}$$

$$\geq \frac{3\gamma}{8} y - \bar{x}^{2}, \tag{B.4}$$

where the second inequality follows from the quadratic growth bound and the third follows from the bound $y - \bar{x} \le 2\delta_{A} \le \frac{3y}{4p}$. Therefore, for all $y \in M \cap B_2\delta_{A}(\bar{x})$, we have

$$f_{M}(\bar{x}) \geq f_{M}(y) + \nabla f_{M}(y), \bar{x} - y + \frac{1}{2} \nabla^{2} f_{M}(y)(\bar{x} - y), (\bar{x} - y) - \frac{\rho}{6} y - \bar{x}^{3}$$

$$\geq f_{M}(y) + \nabla f_{M}(y), \bar{x} - y + \frac{1}{2} \nabla^{2} f_{M}(\bar{x})(\bar{x} - y), (\bar{x} - y) - \frac{2\rho}{3} y - \bar{x}^{3}$$

$$\geq f_{M}(y) + \nabla f_{M}(y), \bar{x} - y + \frac{3\gamma}{8} y - \bar{x}^{2} - \frac{2\rho}{3} y - \bar{x}^{3}$$

$$\geq f_{M}(y) + \nabla f_{M}(y), \bar{x} - y ,$$

where the first and second inequalities follow by Lipschitz continuity of $\nabla^2 fM$; the third inequality follows from (B.4); and the fourth inequality follows from the bound $y - \bar{x} \le 2\delta_{\rm A} \le \frac{9Y}{16D}$. This completes the proof.

B.3 Item 4: Consequences of Strong (a)-Regularity

Fix $x \in B\delta_A(\overline{x})$ and $\sigma \le \delta_A$. Recall that $y := PM(x) \in B_2\delta_A(\overline{x})$ since $x \in B\delta_A(\overline{x})$. Fix $g \in \partial \sigma f(x)$. By definition of $\partial \sigma f(x)$, there exists a family of coefficients $\lambda_i \in [0, 1]$, points $x_i \in \overline{B}\sigma(x) \subseteq \overline{B}\delta(\overline{x})$, and subgradients $g_i \in \partial f(x_i)$ indexed by a finite set $i \in I$ such that $\lim_{i \in I} \lambda_i = 1$ and $g = \lim_{i \in I} \lambda_i g_i$. Therefore, by averaging the strong f(x) bound f(x) over f(x) we find that

$$\begin{split} P_{T_{M}}\left(y\right)&\left(g-\nabla_{M}\ f(y)\right)\leq &\lambda_{i}\ P_{T_{M}}\left(y\right)\left(g_{i}-\nabla_{M}\ f(y)\right)\\ &\leq &\lambda_{i}C_{(a)}\ x_{i}-y.\\ &\stackrel{i\in I}{\leq} &C_{(a)}(\mathrm{dist}(x,M)+\sigma). \end{split}$$

Since *g* was arbitrary, it follows that for all $x \in B\delta_A(\vec{x})$ and $\sigma \le \delta_A$, we have

$$\sup_{g \in \partial \sigma f(x)} P_{TM(y)}(g - \nabla_M f(y)) \le C_{(a)}(\operatorname{dist}(x, M) + \sigma). \tag{B.5}$$

Now we apply this bound to establish the two remaining inequalities.

Indeed, first observe that for all $x \in B\delta_A(\vec{x})$ and $\sigma \le \delta_A$, we have

$$\sup_{g \in \partial \sigma f(x)} P_{TM}(y)g \leq \nabla M f(y) + C_{(a)}(\operatorname{dist}(x, M) + \sigma)$$

$$\leq \beta y - \overline{x} + C_{(a)}(\operatorname{dist}(x, M) + \sigma),$$

where the first inequality follows from (B.5) and the second inequality follows from Item 3. This proves the first claimed bound. Second, observe that for all $x \in B\delta_A(\vec{x})$ and $\sigma \leq \delta_A$, we have

$$\sup_{g,g\in\partial\sigma\,f(x)}P_{T_{M}\ (y)}(g-g)\leq \sup_{g\in\partial\sigma\,f(x)}P_{T_{M}\ (y)}(g-\nabla\,M\ f(y))$$

$$+\sup_{g\in\partial\sigma\,f(x)}P_{T_{M}\ (y)}(g-\nabla\,M\ f(y))$$

$$\leq 2C_{(g)}(\operatorname{dist}(x,M)+\sigma).$$

where the second inequality follows from (B.5). This completes the proof.

B.4 Item 5: Aiming Inequality

Consider a point $x \in B_{\delta_A}(x)$, let $\kappa = 2\mu$, and define

$$\hat{x} \in \underset{x \in \overline{B}_2 \delta_{\Delta}(\overline{x})}{\operatorname{argmin}} f(x) + \kappa x - x$$
.

We claim that $\hat{x} \in M \cap B_2\delta_A(\vec{x})$. Indeed, first note that by definition of \hat{x} and the inclusion $\hat{x} \in \overline{B}_2\delta_A(\vec{x})$, we have

$$\hat{x} - x \le \frac{f(\bar{x}) - f(\hat{x})}{\kappa} + \bar{x} - x \le \bar{x} - x < \delta_{A'}$$

where the second inequality follows since \bar{x} is a minimizer of f on $B_2\delta_A(\bar{x})$, a consequence of quadratic growth. Thus, by the triangle inequality, we have $\hat{x} \in B_2\delta_A(\bar{x})$. By Fermat's rule, we, therefore, have the inclusion:

$$0 \in \partial (\ f + \kappa \cdot - \ x)(\widehat{\ } x) \subseteq \partial \ f(\widehat{x}) + \kappa \, \overline{B} \cdot$$

If $\hat{x} \notin M$, then $\operatorname{dist}(0, \partial f(\hat{x})) > K$, contradicting the above inclusion. Therefore, we have $\hat{x} \in M \cap B_2 \delta_{\Delta}(\bar{x})$.

Turning to the aiming inequality, apply the $(b \le)$ -regularity bound (B.2) to \hat{x} :

$$f(\hat{x}) \ge f(x) + v, \hat{x} - x - \varepsilon \quad x - \hat{x} \ge f(\hat{x}) + v, \hat{x} - x + (\kappa - \varepsilon) \quad x - \hat{x},$$

where we define $\varepsilon := \mu / 2$. Consequently, we have

$$V, x - PM(x) \ge (K - \varepsilon) \quad x - \hat{x} + V, \hat{x} - PM(x)$$
 for all $V \in \partial f(x)$. (B.6)

We now bound the term v, $\hat{x} - PM(x)$: By the conclusion of Item 2, we have

$$PM(\hat{x}) - PM(x) - P_{TM(PM(x))}(\hat{x} - x)$$

 $\leq CM(x - \hat{x}^2 + \text{dist}^2(x, M)) \leq 2CM(x - \hat{x}^2, M)$

where the second inequality follows since $\hat{x} \in M$. Thus, we have

$$| v, \hat{x} - P_{M}(x) | \leq | v, P_{T_{M}(P_{M}(x))}(\hat{x} - x) | + 2C_{M}v \quad x - \hat{x}^{2}$$

$$\leq P_{T_{M}(P_{M}(x))}v^{\hat{x}} - x + 2C_{M}L \quad x - \hat{x}^{2}$$

$$\leq (C_{(a)}\operatorname{dist}(x, M) + \beta \quad P_{M}(x) - \bar{x})^{\hat{x}} - x + 2C_{M}L \quad x - \hat{x}^{2}$$

$$\leq (C_{(a)}\delta_{A} + 2\beta\delta_{A} + 2C_{M}L\delta_{A})^{\hat{x}} - x$$

$$\leq \varepsilon^{\hat{x}} - x.$$

where the second inequality follows from Item 4 and the third inequality follows from the inclusion $P_M(x) \in B_2\delta_A(\vec{x})$. Therefore, plugging this bound into (B.6), we arrive at

$$V, x - P_M(x) \ge (\kappa - 2\varepsilon) x - \hat{x} \ge \mu \operatorname{dist}(x, M),$$

as desired.

B.5 Item 6: Bounding Subgradients

Fix $x \in B\delta_A(\overrightarrow{x})$, $\sigma \le \delta_A$, and $g \in \partial \sigma f(x)$. By definition of $\partial \sigma f(x)$, there exists a family of coefficients $\lambda_i \in [0, 1]$, points $x_i \in \overline{B}\sigma(x) \subseteq \overline{B}\delta(\overrightarrow{x})$, and subgradients $g_i \in \partial f(x_i)$ indexed by a finite set $i \in I$ such that $f(x_i) = 1$ and $f(x_i) = 1$ and $f(x_i) = 1$. Recall that by Lipschitz continuity of $f(x_i) = 1$ on $f(x_i) = 1$. Therefore,

$$g \leq \lambda_i \ g_i \leq L,$$

as desired.

B.6 Item 7: Bounding the Function Gap

Fix a point $x \in B_{\delta_A}(\vec{x})$ and recall that $P_M(x) \in B_2\delta_A(\vec{x})$. Then by Lipschitz continuity of f on $B_\delta(\vec{x})$, we have

$$f(x) - f(P_M(\bar{x})) \le L \operatorname{dist}(x, M).$$

Next, arguing as in the proof of Item 3, we find that ∇ fM (\overline{x}) = 0. Thus, since ∇ fM is β -Lipschitz on $B\delta(\overline{x})$, we have

$$f(PM(x)) - f(\bar{x}) = fM(PM(x)) - f(\bar{x}) \le \nabla fM(\bar{x}), PM(x) - \bar{x} + \frac{\beta}{2} PM(x) - \bar{x}^2 = \frac{\beta}{2} PM(x) - \bar{x}^2.$$

By putting both bounds together, we have

$$f(x) - f(\overline{x}) = f(x) - f(P_M(x)) + f(P_M(x)) - f(\overline{x})$$

 $\leq L \text{dist}(x, M) + \frac{\beta}{2} P_M(x) - \overline{x}^2,$

as desired.

C Proof of Corollary 2.5

We begin with the following known Lemma, which immediately follows from [29, Proposition 2.8]

Lemma C.1 Let $\underline{f}: \mathbb{R}^d \to \mathbb{R}$ be a locally Lipschitz function. Suppose that there exists sequences $x_k \to \overline{x}$, $\tau_k \to 0$, and $g_k \in \partial \tau_k f(x_k)$ with $g_k \to 0$. Then \overline{x} is a Clarke critical point.

Now we turn to the proof of the Corollary. Since *f* has bounded initial sublevel set, the following widened sublevel set is bounded:

$$S := \{ x + u : f(x) \le f(x_0) \text{ and } u \in \overline{B}(x) \}.$$

Thus, there exists $L \ge 0$ such that f is L-Lipschitz on S. In addition, ∂f is uniformly bounded by L on int S.

We begin with a claim.

<u>Claim</u>: Fix i > 0 and define $\tau_i := 2^{-i}$. Let $s_k := \max\{\sqrt{\frac{g_k}{128L}}, c_0 g_0\}$ be the trust region parameter used in Algorithm 3 and define $s_{i,k} := \frac{1}{128L}, \tau_i$. Then with probability one, the event

$$E_k^{(i)} = \text{dist}(0, \, \partial_{\tau_i} f(x_k)) > \sum_{i,k \text{ and } f(x_{k+1})} f(x_k) - \frac{\tau_i \operatorname{dist}(0, \, \partial_{\tau_i} f(x_k))}{8}$$

cannot happen infinitely often, i.e.,

$$P \cap_{T=1}^{\infty} \cup_{k=T}^{\infty} E_k^{(i)} = 0$$

Proof We prove that $P(E_k^{(i)})$ is summable in k. Indeed, first, note that $P(E_k^{(i)}) = 0$ when $P(\text{dist}(0, \partial_{\tau_i} f(x_k))) > 0$. On the other hand, suppose

 $P(\operatorname{dist}(0, \partial_{\tau_i} f(x_k)) > i,k) > 0$. Now we upper bound $P(E_k^{(i)})$ for all G_k satisfying $G_k \ge i$. For such $G := G_k$, the radius $\tau_i = \sigma_{G-i}$ is among those considered in Algorithm 3. Moreover, since $s_k \le L$ (recall $x_k \in \operatorname{int}(S)$), the radius satisfies the trust region constraint: $\sigma_{G-i} = \tau_i \le i,k's_k \le \operatorname{dist}(0,\partial_{\sigma_{G-i}} f(x))/s_k$. Therefore, if NDescent terminates with descent at the (G-i)-th level in Algorithm 3, it follows that

$$f(x_{k+1}) > f(x_k) - \frac{\tau_i \operatorname{dist}(0, \, \partial_{\tau_i} f(x_k))}{8}.$$

We estimate the probability of this success with Lemma 2.3: there exist \mathcal{E} 0 depending on i,k and for all $k \geq i$, we have

$$P(E_k^{(i)}) \le P \quad f(x_{k+1}) > f(x_k) - \frac{\tau_i \operatorname{dist}(0, \, \partial_{\tau_i} f(x_k))}{8} || \operatorname{dist}(0, \, \partial_{\tau_i} f(x_k)) > |_{i,k}$$

$$\le \exp(-Ck).$$

Therefore, $P(E_k^{(i)})$ is summable in k. The result then follows from Borel–Cantelli lemma.

By the claim and a union bound, we know that with probability one, for any fixed i, $E_k^{(i)}$ cannot happen infinitely often. Now, suppose that a subsequence $\{x_{k_l}\}$ (where $k_l \ge l$ is strictly increasing in l) converges to a point x. We note that the sequence $\{f(x_k)\}$ is bounded below: Indeed, since x_{k_l} converges and f is continuous, it follows $\{f(x_{k_l})\}$ is bounded below by a constant $c \in R$. Consequently, since $\{f(x_k)\}$ is nonincreasing and $k_l \ge l$, it follows that $c \le f(x_{k_l}) \le f(x_l)$ and for every $l \ge 0$, as desired. As a result, the following inequalities cannot be valid simultaneously infinitely often:

$$\operatorname{dist}(0, \, \partial_{\tau_i} f(x_{k_l})) > \prod_{i,k \text{ and } f} (x_{k_l+1}) \le f(x_{k_l}) - \frac{\tau_i \operatorname{dist}(0, \, \partial_{\tau_i} f(x_{k_l}))}{8}.$$

Therefore, dist $(0, \partial_{\tau_i} f(x_{k_l})) > i,k$ cannot happen infinitely often. Consequently, we can find a sequence of increasing indices j_i such that

$$\operatorname{dist}(0, \partial_{\tau_i} f(x_{j_i})) \leq i,k \quad \text{and } x_{j_i} \to x$$

Since $i,k \to 0$ as $k \to \infty$, Lemma C.1, shows that \overline{x} is Clarke critical.

D Proof of Lemma 5.7

We begin with preliminary notation and bounds. We fix $x \in B\delta_{Grid}(x)$ and subgradient $g \in \partial_{\sigma} f(x) \setminus \{0\}$. We define y := PM(x), T := TM(y), and N := NM(y). We have the following two bounds: First, we have

$$(\mu + L)C_M \left(D_1 \mathrm{dist}(x,M) + \sigma\right) \leq (\mu + L)C_M \delta_{\mathrm{Grid}}(D_1 + 1)$$

$$= \frac{\mu}{8} C_M (D_1^{-1} + 1) \delta_{Grid} \le \frac{\mu}{8}.$$
 (D.1)

Second, we have

$$C_{(a)}(\operatorname{dist}(x,M) + \sigma) + \beta \quad y - \bar{x} \le 2C_{(a)}\delta_{\operatorname{Grid}} + 2\beta\delta_{\operatorname{Grid}} \le \frac{\mu}{4}. \tag{D.2}$$

We now turn to the proof.

By Lemma 5.3 (which is applicable since $x \in B_{\delta_A/2}(\vec{x})$ and $\sigma \le \delta_{Grid} \le \delta_A/2$), we have

$$\hat{g}, \, \sigma \frac{P_N g}{g}^! \leq -\sigma \mu \, \frac{P_N g}{g} + (\mu + \ L) \mathrm{dist}(x, \, M) + (\mu + \ L) C_M \, (\mathrm{dist}^2(x, \, M) + \sigma^2).$$

Rearranging, we find that

$$\begin{split} P_N \hat{g} \cdot g &\leq -\mu &\quad P_N g + \quad \frac{(\mu + \ L) \ g \ \operatorname{dist}(x, \ M)}{\sigma} + \frac{(\mu + \ L) \ g \ C_M \ (\operatorname{dist}^2(x, \ M) + \sigma^2)}{\sigma} \\ &\leq -\mu &\quad P_N g + \quad \frac{\mu}{8} \ g + (\mu + \ L) C_M \ (D_1 \operatorname{dist}(x, \ M) + \sigma) \cdot \ g \\ &\leq -\mu &\quad P_N g + \quad \frac{\mu}{4} \ g, \end{split}$$

where the second inequality follows from the assumption $D_1^{-1} \operatorname{dist}(x, M) \le \sigma$ and the third follows from (D.1). Now observe that

$$P_T\hat{g}, g \le P_T\hat{g} \quad g \le (C_{(a)}(\operatorname{dist}(x, M) + \sigma) + \beta \quad y - \bar{x}) \cdot g \le \frac{\mu}{4} g,$$

where second inequality follows from (3.4) and the third inequality follows from (D.2). Therefore,

$$\hat{g}, g = P_N \hat{g}, g + P_T \hat{g}, g \leq -\mu \quad P_N g + \frac{\mu}{2} g \leq -\frac{\mu}{2} g + \mu \quad P_T(g),$$

as desired.

E Proof that≤ L

Lemma E.1 We have that $\mu \leq L$.

Proof Indeed,

$$\mu = \frac{1}{4} \lim_{\substack{x \to -\infty \\ x \to -x}} \inf_{x} \operatorname{dist}(0, \, \partial f(x)) \leq \lim_{x \to -\infty} \sup_{x} \operatorname{dist}(0, \, \partial f(x)) \leq L.$$

by Proposition 3.5.

F Proof of Lemma 6.4

We fix a > 0. Note that the claimed inclusion is a consequence of the following bound:

$$f(x) - f(\overline{x}) \ge \frac{\gamma}{2} \min{\{\delta_A, x - \overline{x}\}} \quad x - \overline{x} \quad \text{for all } x \in \mathbb{R}^d.$$
 (F.1)

Here we provide a proof for completeness.

To that end, we remind the reader that Assumption A is in force. Consequently, by Item 1 of Proposition 3.5, we have:

$$f(x) - f(\bar{x}) \ge \frac{\gamma}{2} |x - \bar{x}|^2$$
 for all $x \in \overline{B}_{\delta_A}(\bar{x})$.

Thus, if $x \in B\delta_A(\vec{x})$, bound (F.1) is immediate. On the other hand, suppose that we have $x \in \mathbb{R}^d \setminus B\delta_A(\vec{x})$. Define the curve $x_t : t \to (1-t)x + t\bar{x}$. Choose $t_0 \in [0, 1]$ such that $x_{t_0} \in \text{bdry } B\delta_A(\vec{x})$. Then by Jensen's inequality, we have

$$(1-t_0) f(x) \ge f(x_{t_0}) - t_0 f(\bar{x}) \ge (1-t_0) f(\bar{x}) + \frac{\gamma}{2} x_{t_0} - \bar{x}^{-2}$$

$$= (1-t_0) f(\bar{x}) + \frac{\gamma (1-t_0)}{2} x - \bar{x} x_{t_0} - \bar{x}.$$

Consequently, since $x_{t_0} - \overline{x} = \delta_A$, we have

$$f(x) - f(\overline{x}) \ge \frac{y \, \delta_A}{2} x - \overline{x} \ge \frac{y}{2} \min\{\delta_A, x - \overline{x}\} x - \overline{x},$$

as desired. This completes the proof.

G Proof of (6.18)

Let us expand the left-hand-side of (6.18):

$$\frac{16L \quad \overline{2 \log(2K_1^{2/p})}}{K_1^{1/2}} \leq \frac{16DL \quad \overline{2 \log(K_1^2)}}{\|\underbrace{K_1^{1/2}}_{=:A} - \|} + \frac{16DL \quad \overline{2 \log(2/p)}}{\|\underbrace{K_1^{1/2}}_{=:B} - \|}.$$

Note that $B \le a/4$ by definition of K_1 . Consequently, the proof will follow if $A \le a/4$. To that end, for any $\alpha \in (0, 1)$, we have

$$A = \frac{16DL}{K_1^{1/2}} \frac{\overline{2 \log(K_1^2)}}{K_1^{1/2}} = \frac{16DL}{K_1^{1/2}} \frac{\overline{2 \log(K_1^{2\alpha})/\alpha}}{K_1^{1/2}} \leq \frac{16DL}{K_1^{(1-\alpha)/2}},$$

Therefore, we have $A \le a/4$ whenever

$$K_1 \geq \inf_{\alpha \in \{0,1\}} \frac{-64DL \quad \frac{\frac{2}{\alpha}}{a} \quad \frac{\frac{2}{(1-\alpha)}}{a}}{a^{\frac{2}{(1-\alpha)}}} = \frac{D^2}{a^2} \inf_{\alpha \in \{0,1\}} \frac{-64L \quad \frac{\frac{2}{\alpha}}{a} \quad \frac{\frac{2}{(1-\alpha)}}{a}}{0 \frac{a}{D} \cdot \frac{1}{(1-\alpha)}} = \frac{D^2}{a^2} b.$$

This lower bound holds by definition of K_1 . Consequently $A \le a/4$. Therefore, the proof is complete.

H Proof of Lemma 6.6

Throughout this section, we use the symbol a b to mean that $a \le \eta b$ for a fixed numerical constant η that is independent of f. In addition, we use the bound on the condition number: $\kappa \ge 1$, since $\mu \le L$; see Lemma E.

Turning to the bound, we wish to upper bound q.

$$q = \max \left\{ \rho, \frac{3\mu^2}{1 - \frac{3\mu^2}{256L^2}}, \frac{1}{2} \right\}$$

First note that

$$1 - \frac{3\mu^2}{1 - \frac{3\mu^2}{256L^2}} \quad \frac{\mu^2}{L^2} \ge \frac{1}{\kappa^2}.$$

Next, we upper bound ρ . To that end, we must bound the constants a_1 and a_2 , which rely on the somewhat involved constants C_4 and C_5 . Thus, we first lower bound C_4 :

$$C_{4} = \min \frac{\beta}{C_{(a)}(1 + \delta_{A})}, \frac{\min \{\mu/\delta_{A}, C_{3}D_{2}/\beta\}}{4(1 + (1 + \delta_{A})C_{M})(\mu + L))}, \frac{1}{2}$$

$$\min \frac{\beta}{C_{(a)}}, \frac{\mu}{L(1 + C_{M})}, \frac{y^{2}\mu}{L^{2}\beta(1 + C_{M})}$$

$$\geq \frac{1}{\kappa^{3}(1 + C_{M})},$$

where we use the bounds $\mu \le L$, C_3 γ^{-2}/L , and D_2 μ^- . Turning to C_5 , we have:

$$\begin{split} C_5 &= \min \ \frac{\beta}{2C_{(a)}}, \frac{C_3D_2}{32C_{(a)}\beta}, C_4, \frac{C_2}{4} \\ & \min \ \frac{\beta}{C_{(a)}}, \frac{\gamma^2\mu}{LC_{(a)}\beta}, \frac{1}{\kappa^3(1+C_M)}, \frac{\gamma}{C_{(a)}} \\ & \geq \frac{1}{\kappa^3(1+C_M)}, \end{split}$$

where we again use C_3 y $^{2}/L$, and D_2 μ . Therefore, we have the lower bound for a_2 :

$$a_2 = \frac{\min \{C_1/L, C_5\}}{2} \quad \min \quad \frac{Y^2}{L^2}, \frac{1}{K^3(1+C_M)} \qquad \frac{1}{K^3(1+C_M)}.$$

In addition, we have the upper bound:

$$a_2 = \frac{\min\{C_1/L, C_5\}}{2} \le C_4/2 \le 1/4.$$

Finally to lower bound a_1 , we have

$$a_1 = \min\{D_1, D_2/L\} \quad \min \quad \frac{\mu}{L}, \frac{\gamma}{L} \quad \frac{1}{\kappa},$$

where we use the bound $D_1 \mu / L$ and $D_2 / L \mu / L$. Now we upper bound ρ by providing a lower bound on $1 - \rho$.

$$1 - \rho = \frac{1}{8} \min \frac{y_{a_2}}{8 \max\{4La_2^2, \beta\}}, \frac{\mu_{a_1}}{4 \max\{2L, \beta/a_2^2\}}$$

$$\min \frac{y}{La_2}, \frac{y_{a_2}}{\beta}, \frac{\mu_{a_1}}{L}, \frac{\mu_{a_1}a_2^2}{\beta}$$

$$\min \frac{y}{L}, \frac{y}{\kappa^3(1 + C_M)\beta}, \frac{\mu}{\kappa_L}, \frac{\mu}{\beta\kappa^7(1 + C_M)^2}$$

$$\frac{1}{\kappa^8(1 + C_M)^2}$$

Putting all these bounds together, we find that:

$$1 - q \quad \min \quad \frac{1}{\kappa^8 (1 + C_M)^2}, \frac{1}{\kappa^2} \geq \frac{1}{\kappa^8 (1 + C_M)^2},$$

as desired.

References

- 1. K. M. Anstreicher and J. Lee. A masked spectral bound for maximum-entropy sampling. In *mODa* 7–Advances in Model-Oriented Design and Analysis, pages 1–12. Springer, 2004.
- 2. H. Attouch and J. Bolte. On the convergence of the proximal algorithm for nonsmooth functions involving analytic features. *Mathematical Programming*, 116:5–16, 2009.
- H. Attouch, J. Bolte, P. Redont, and A. Soubeyran. Proximal alternating minimization and projection methods for nonconvex problems: An approach based on the Kurdyka-Łojasiewicz inequality.
 Mathematics of operations research, 35(2):438–457, 2010.

- H. Attouch, J. Bolte, and B. F. Svaiter. Convergence of descent methods for semi-algebraic and tame problems: proximal algorithms, forward–backward splitting, and regularized gauss–seidel methods. *Mathematical Programming*, 137(1):91–129, 2013.
- 5. P. Bianchi, W. Hachem, and S. Schechtman. Stochastic subgradient descent escapes active strict saddles. *arXiv preprint* arXiv:2108.02072, 2021.
- J. Bolte, A. Daniilidis, and A. Lewis. The Łojasiewicz inequality for nonsmooth subanalytic functions with applications to subgradient dynamical systems. SIAM J. Optim., 17(4):1205–1223 (electronic), 2006.
- J. Bolte, A. Daniilidis, A. Lewis, and M. Shiota. Clarke subgradients of stratifiable functions. SIAM Journal on Optimization. 18(2):556–572, 2007.
- 8. J. Bolte, T. P. Nguyen, J. Peypouquet, and B. W. Suter. From error bounds to the complexity of first-order descent methods for convex functions. *Mathematical Programming*, 165:471–507, 2017.
- 9. J. Bolte, S. Sabach, and M. Teboulle. Proximal alternating linearized minimization for nonconvex and nonsmooth problems. *Mathematical Programming*, 146(1):459–494, 2014.
- 10. N. Boumal. An introduction to optimization on smooth manifolds. Available online, Aug, 2020.
- 11. S. Bubeck et al. Convex optimization: Algorithms and complexity. Foundations and Trends® in Machine Learning, 8(3-4):231–357, 2015.
- 12. S. Burer and J. Lee. Solving maximum-entropy sampling problems using factored masks. *Mathematical Programming*, 109(2):263–281, 2007.
- 13. J. V. Burke. Descent methods for composite nondifferentiable optimization problems. *Mathematical Programming*, 33(3):260–279, 1985.
- V. Charisopoulos, Y. Chen, D. Davis, M. Díaz, L. Ding, and D. Drusvyatskiy. Low-rank matrix recovery with composite optimization: good conditioning and rapid convergence. *Foundations of Computational Mathematics*, 21(6):1505–1593, 2021.
- 15. Y. Chen, Y. Chi, and A. J. Goldsmith. Exact and stable covariance estimation from quadratic sampling via convex programming. *IEEE Transactions on Information Theory*, 61(7):4034–4059, 2015.
- 16. F. H. Clarke, Y. S. Ledyaev, R. J. Stern, and P. R. Wolenski. *Nonsmooth analysis and control theory*, volume 178. Springer Science & Business Media, 2008.
- 17. D. Davis and D. Drusvyatskiy. Subgradient methods under weak convexity and tame geometry. *SIAG/OPT Views and News*, 28(1):1–10, 2020.
- D. Davis, D. Drusvyatskiy, and V. Charisopoulos. Stochastic algorithms with geometric step decay converge linearly on sharp functions. arXiv preprint arXiv:1907.09547, 2019.
- 19. D. Davis, D. Drusvyatskiy, and L. Jiang. Active manifolds, stratifications, and convergence to local minima in nonsmooth optimization. *arXiv preprint* arXiv:2108.11832, 2023.
- 20. D. Davis, D. Drusvyatskiy, S. Kakade, and J. D. Lee. Stochastic subgradient method converges on tame functions. *Foundations of computational mathematics*, 20(1):119–154, 2020.
- 21. D. Davis, D. Drusvyatskiy, Y. T. Lee, S. Padmanabhan, and G. Ye. A gradient sampling method with complexity guarantees for general lipschitz functions. *arXiv preprint* arXiv:2112.06969, 2022.
- 22. D. Davis, D. Drusvyatskiy, K. J. MacPhee, and C. Paquette. Subgradient methods for sharp weakly convex functions. *Journal of Optimization Theory and Applications*, 179:962–982, 2018.
- D. Drusvyatskiy, A. D. Ioffe, and A. S. Lewis. Generic minimizing behavior in semialgebraic optimization. SIAM Journal on Optimization, 26(1):513–534, 2016.
- D. Drusvyatskiy, A. D. Ioffe, and A. S. Lewis. Nonsmooth optimization using taylor-like models: error bounds, convergence, and termination criteria. *Mathematical Programming*, 185:357–383, 2021.
- 25. D. Drusvyatskiy and A. Lewis. Error bounds, quadratic growth, and linear convergence of proximal methods. *To appear in Math. Oper. Res.*, arXiv:1602.06661, 2016.
- 26. D. Drusvyatskiy and A. S. Lewis. Optimality, identifiability, and sensitivity. *Mathematical Programming*, 147(1):467–498, 2014. Citations refer to long version arXiv:1207.6628.
- P. J. Enright and B. A. Conway. Discrete approximations to optimal trajectories using direct transcription and nonlinear programming. *Journal of Guidance, Control, and Dynamics*, 15(4):994–1002, 1992.
- 28. R. Fletcher. A model algorithm for composite nondifferentiable optimization problems. In *Nondifferential and Variational Techniques in Optimization*, pages 67–76. Springer, 1982.
- 29. A. Goldstein. Optimization of lipschitz continuous functions. *Mathematical Programming*, 13(1):14–22, 1977.
- 30. X. Han and A. S. Lewis. Survey descent: A multipoint generalization of gradient descent for nonsmooth optimization. *arXiv preprint* arXiv:2111.15645, 2021.

- 31. W. L. Hare and A. S. Lewis. Identifying active constraints via partial smoothness and prox-regularity. *Journal of Convex Analysis*, 11(2):251–266, 2004.
- 32. A. Ioffe, An invitation to tame optimization. SIAM J. Optim., 19(4):1894–1917, 2009.
- 33. P. R. Johnstone and P. Moulin. Faster subgradient methods for functions with hölderian growth. *Mathematical Programming*, 180(1-2):417–450, 2020.
- 34. J. M. Lee. Smooth manifolds. In Introduction to Smooth Manifolds, pages 1–31. Springer, 2013.
- 35. C. Lemarechal. An extension of davidon methods to non differentiable problems. In *Nondifferentiable optimization*, pages 95–109. Springer, 1975.
- 36. C. Lemaréchal, F. Oustry, and C. Sagastizábal. The *U*-lagrangian of a convex function. *Transactions of the American mathematical Society*, 352(2):711–729, 2000.
- 37. A. Lewis and C. Wylie. A simple newton method for local nonsmooth optimization. *arXiv* preprint arXiv:1907.11742, 2019.
- 38. A. S. Lewis. Active sets, nonsmoothness, and sensitivity. SIAM Journal on Optimization, 13(3):702–725, 2002.
- 39. A. S. Lewis and M. L. Overton. Nonsmooth optimization via quasi-newton methods. *Mathematical Programming*, 141(1):135–163, 2013.
- 40. R. Mifflin. Semismooth and semiconvex functions in constrained optimization. *SIAM J. Control Optim.*, 15(6):959–972, 1977.
- 41. R. Mifflin and C. Sagastizábal. AVU-algorithm for convex minimization. *Mathematical Programming*, 104(2):583–608, 2005.
- 42. A. Nemirovsky and D. Yudin. *Problem complexity and method efficiency in optimization*. A Wiley-Interscience Publication. John Wiley & Sons, Inc., New York, 1983. Translated from the Russian and with a preface by E. R. Dawson, Wiley-Interscience Series in Discrete Mathematics.
- 43. Y. Nesterov. *Introductory lectures on convex optimization*, volume 87 of *Applied Optimization*. Kluwer Academic Publishers, Boston, MA, 2004. A basic course.
- 44. W. D. Oliveira and C. Sagastizábal. Bundle methods in the xxist century: A bird's-eye view. *Pesquisa Operacional*, 34:647–670, 2014.
- A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein,
 L. Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- 46. R. Poliquin and R. Rockafellar. Prox-regular functions in variational analysis. *Transactions of the American Mathematical Society*, 348(5):1805–1838, 1996.
- 47. B. T. Polyak. Minimization of unsmooth functionals. *USSR Computational Mathematics and Mathematical Physics*, 9(3):14–29, 1969.
- 48. R. Rockafellar and R.-B. Wets. *Variational Analysis*. Grundlehren der mathematischen Wissenschaften, Vol 317, Springer, Berlin, 1998.
- 49. A. Shapiro. On a class of nonsmooth composite functions. *Mathematics of Operations Research*, 28(4):677–692, 2003.
- 50. P. Wolfe. A method of conjugate subgradients for minimizing nondifferentiable functions. In *Nondifferentiable optimization*, pages 145–173. Springer, 1975.
- 51. S. J. Wright. Convergence of an inexact algorithm for composite nonsmooth optimization. *IMA journal of numerical analysis*, 10(3):299–321, 1990.
- 52. S. J. Wright. Identifiable surfaces in constrained optimization. SIAM Journal on Control and Optimization, 31(4):1063–1079, 1993.
- 53. Y. Xu and W. Yin. A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion. *SIAM Journal on imaging sciences*, 6(3):1758–1789, 2013.
- 54. Y.-X. Yuan. On the superlinear convergence of a trust region algorithm for nonsmooth optimization. *Mathematical Programming*, 31(3):269–285, 1985.
- J. Zhang, H. Lin, S. Jegelka, S. Sra, and A. Jadbabaie. Complexity of finding stationary points of nonconvex nonsmooth functions. *Proceedings of Machine Learning Research*, pages 11173–11182, 2020.

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.