

# Gradient Sampling Methods for Nonsmooth Optimization

J.V. Burke\*   F.E. Curtis†   A.S. Lewis‡   M.L. Overton§   L.E.A. Simões¶

May 1, 2018

Submitted to: *Special Methods for Nonsmooth Optimization*, Springer, 2018  
A. Bagirov, M. Gaudioso, N. Karmitsa and M. Mäkelä, eds.

## Abstract

This paper reviews the gradient sampling methodology for solving nonsmooth, nonconvex optimization problems. An intuitively straightforward gradient sampling algorithm is stated and its convergence properties are summarized. Throughout this discussion, we emphasize the simplicity of gradient sampling as an extension of the steepest descent method for minimizing smooth objectives. We then provide overviews of various enhancements that have been proposed to improve practical performance, as well as of several extensions that have been made in the literature, such as to solve constrained problems. The paper also includes clarification of certain technical aspects of the analysis of gradient sampling algorithms, most notably related to the assumptions one needs to make about the set of points at which the objective is continuously differentiable. Finally, we discuss possible future research directions.

## 1 Introduction

The Gradient Sampling (GS) algorithm is a conceptually simple descent method for solving nonsmooth, nonconvex optimization problems, yet it is one that possesses a solid theoretical foundation and has been employed to substantial success in a wide variety of applications. Since the appearance of the fundamental algorithm and its analysis little over a dozen years ago, GS has matured into a comprehensive methodology. Various enhancements have been proposed that make it a competitive approach in many nonsmooth optimization contexts, and it has been extended in various interesting ways, such as for nonsmooth optimization on manifolds and for constrained problems. The purpose of this work is to provide background and motivation for the development of the GS method, discuss its theoretical guarantees, and provide an overview of the enhancements and extensions that have been the subject of research over recent years.

The underlying philosophy of GS is that virtually any nonsmooth objective function of interest is differentiable almost everywhere; in particular, this is true if the objective  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is either locally Lipschitz continuous or semialgebraic. In such cases, when  $f$  is evaluated at a randomly generated point  $x \in \mathbb{R}^n$ , it is differentiable there with probability one. This means that an algorithm can rely on an ability to obtain the objective function value  $f(x)$  and gradient  $\nabla f(x)$ , as when  $f$  is smooth, rather than require an oracle to compute a subgradient. In most interesting settings,  $f$  is *not* differentiable at its local minimizers, but,

---

\*Department of Mathematics, University of Washington, Seattle, WA. [jvburke01@gmail.com](mailto:jvburke01@gmail.com). Supported in part by the U.S. National Science Foundation grant DMS-1514559.

†Department of Industrial and Systems Engineering, Lehigh University, Bethlehem, PA. [frank.e.curtis@gmail.com](mailto:frank.e.curtis@gmail.com). Supported in part by the U.S. Department of Energy grant DE-SC0010615.

‡School of Operations Research and Information Engineering, Cornell University, Ithaca, NY. [adrian.lewis@cornell.edu](mailto:adrian.lewis@cornell.edu). Supported in part by the U.S. National Science Foundation grant DMS-1613996.

§Courant Institute of Mathematical Sciences, New York University. [mo1@nyu.edu](mailto:mo1@nyu.edu). Supported in part by the U.S. National Science Foundation grant DMS-1620083.

¶Department of Applied Mathematics, University of Campinas, Brazil. [simoes.lea@gmail.com](mailto:simoes.lea@gmail.com). Supported in part by the São Paulo Research Foundation (FAPESP), Brazil, under grants 2016/22989-2 and 2017/07265-0.

under reasonable assumptions, the carefully crafted mechanisms of the GS algorithm generate a sequence of iterates—at which  $f$  is differentiable—converging to stationarity.

At the heart of GS is a stabilized steepest descent approach. When  $f$  is differentiable at  $x$ , the negative gradient  $-\nabla f(x)$  is, of course, the traditional steepest descent direction for  $f$  at  $x$  in the 2-norm in that

$$-\frac{\nabla f(x)}{\|\nabla f(x)\|_2} = \arg \min_{\|d\|_2 \leq 1} \nabla f(x)^T d. \quad (1.1)$$

However, when  $x$  is near a point where  $f$  is not differentiable, it may be necessary to take a very short step along  $-\nabla f(x)$  to obtain decrease in  $f$ . It is for this reason that the traditional steepest descent method may converge to nonstationary points when  $f$  is nonsmooth.<sup>1</sup> The GS algorithm stabilizes the choice of the search direction to avoid this issue. In each iteration, a descent direction from the current iterate  $x^k$  is obtained by supplementing the information provided by  $\nabla f(x^k)$  with gradients evaluated at randomly generated points  $\{x^{k,1}, \dots, x^{k,m}\} \subset \mathbb{B}(x^k, \epsilon_k) := \{x \in \mathbb{R}^n : \|x - x^k\|_2 \leq \epsilon_k\}$ , which are *near*  $x^k$ , and then computing the minimum-norm vector  $g^k$  in the convex hull of these gradients. This choice can be motivated by the goal that, with  $\bar{\partial}_\epsilon f(x)$  denoting the Clarke  $\epsilon$ -subdifferential of  $f$  at  $x$  (see §3),

$$-\frac{g^k}{\|g^k\|_2} \approx \arg \min_{\|d\|_2 \leq 1} \max_{g \in \bar{\partial}_\epsilon f(x)} g^T d; \quad (1.2)$$

i.e.,  $-g^k$  can essentially be viewed as a steepest descent direction for  $f$  from  $x^k$  in a more “robust” sense. A line search is then used to find a positive stepsize  $t_k$  yielding decrease in  $f$ , i.e.,  $f(x^k - t_k g^k) < f(x^k)$ . The sampling radius  $\epsilon_k$  that determines the meaning of “near  $x^k$ ” may either be fixed or adjusted dynamically.

A specific instance of the GS algorithm is presented in §2. Its convergence guarantees are summarized in §3. We then present various enhancements and extensions of the approach in §4 and §5, respectively, followed by a discussion of some successful applications of the GS methodology in §6. Throughout this work, our goal is to emphasize the *simplicity* of the fundamental GS strategy. We believe that this, in addition to its strong convergence properties for locally Lipschitz optimization, makes it an attractive choice when attempting to solve difficult types of nonsmooth optimization problems.

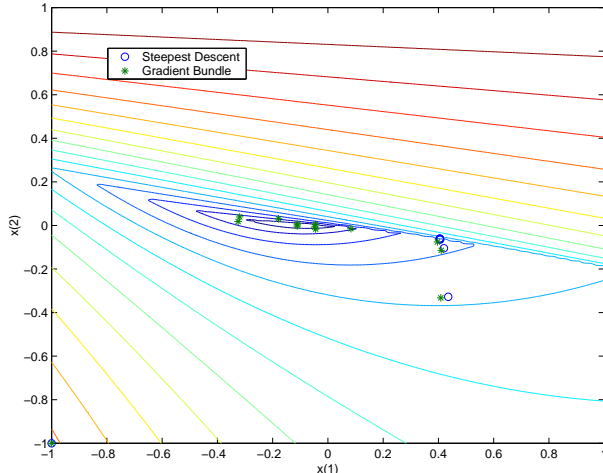
Although the first convergence analysis of a GS algorithm was given by Burke, Lewis, and Overton in [BLO05], an earlier version of the method was presented by these authors in [BLO02b]. That algorithm, originally called a “gradient bundle” method, was applied to a function that was not only nonconvex and nonsmooth, but also non-locally-Lipschitz, namely, the spectral abscissa—i.e., the largest of the real parts of the eigenvalues—of a linear matrix function  $A$  mapping a parameter vector  $x$  to the space of nonsymmetric square matrices. The spectral abscissa is not locally Lipschitz at a matrix  $\bar{X}$  when an eigenvalue of  $\bar{X}$  with largest real part has multiplicity two or more [BO01], but it is semialgebraic and, hence, differentiable almost everywhere, so a GS algorithm was applicable. The method was surprisingly effective. As anticipated, in most cases the apparent local minimizers that were approximated had the property that the eigenvalues of  $A$  with largest real part had multiplicity two or more. An illustration that appeared in [BLO02b] is reproduced in Figure 1; the extremely “steep” contours of the objective function indicate its non-Lipschitzness. Obtaining theoretical results for a GS algorithm applied to non-locally-Lipschitz problems seems very challenging; we discuss this issue further in §3.3, after describing the substantial body of theory that has been developed for the locally Lipschitz case in §3.1 and §3.2.

## 2 Algorithm GS

We now state a specific variant of the GS algorithm. We start by assuming only that the objective function  $f$  is locally Lipschitz over  $\mathbb{R}^n$ , which implies, by Rademacher’s theorem [Cla83], that  $f$  is differentiable almost everywhere. As previously mentioned, at the heart of the algorithm is the computation of a descent direction

---

<sup>1</sup>Although this fact has been known for decades, most of the examples that appear in the literature are rather artificial since they were designed with exact line searches in mind. Analyses showing that the steepest descent method with inexact line searches converges to nonstationary points of some simple convex nonsmooth functions have appeared recently in [AO17, GL17].



**Figure 1** – Contours of the spectral abscissa of an affine matrix family given in [BLO02b]. Iterates of the ordinary gradient (“steepest descent”) method with a line search are shown (small circles) along with those of the gradient sampling (“gradient bundle”) algorithm (asterisks). Both start at  $(-1, -1)$ .

by finding the minimum norm element of the convex hull of gradients obtained about each iterate. The remaining procedures relate to the line search to obtain decrease in  $f$  and the selection of a subsequent iterate so that  $f$  is differentiable at  $\{x^k\}$ .

While the essence of the methods from [BLO02b] and [BLO05] remains intact, Algorithm **GS** differs in subtle yet important ways from the methods presented in these papers, as we now explain.

1. Algorithm **GS** incorporates a key modification proposed by Kiwiel in [Kiw07, Alg. 2.1], namely, the second inequality in (2.2); the version in [BLO05] used  $\epsilon_k$  instead of  $\min\{t_k, \epsilon_k\}$ . As Kiwiel explained, this minor change allowed him to drop the assumption in [BLO05] that the level set  $\{x : f(x) \leq f(x^0)\}$  is compact, strengthening the convergence results for the algorithm.
2. A second change suggested in [Kiw07, Alg. 2.1] is the introduction of the termination tolerances  $\nu_{\text{opt}}$  and  $\epsilon_{\text{opt}}$ . These were used in the computational experiments in [BLO05], but not in the algorithm statement or analysis. Note that if  $\epsilon_{\text{opt}}$  is set to zero, then Algorithm **GS** never terminates since  $\epsilon_k$  can never be zero, though it may happen that one obtains  $\|g^k\|_2 = 0$ .
3. A third change, also suggested by Kiwiel, is the usage of the nonnormalized search direction  $-g^k$  (originally used in [BLO02b]) instead of the normalized search direction  $-g^k/\|g^k\|_2$  (used in [BLO05]). The resulting inequalities in (2.1) and (2.2) are taken from [Kiw07, Sec. 4.1]. This choice does not affect the main conclusions of the convergence theory as in both cases it is established [BLO05, Kiw07] that the stepsize  $t_k$  can be determined by a finite process. However, since Theorem 3.1 below shows that a subsequence of  $\{g^k\}$  converges to zero under reasonable conditions, one expects that fewer function evaluations should be required by the line search asymptotically when using the nonnormalized search direction, whereas using the normalized direction may result in the number of function evaluations growing arbitrarily large [Kiw07, Sec. 4.1]. Furthermore, our practical experience is consistent with this viewpoint.
4. Another aspect of Algorithm **GS** that is different in both [BLO05] and [Kiw07] concerns the randomization procedure in Step 2. In the variants given in those papers, it was stated that the algorithm terminates if  $f$  is not *continuously* differentiable at the randomly sampled points  $\{x^{k,1}, \dots, x^{k,m}\}$ . In the theorem stated in the next section, we require only that  $f$  is differentiable at the sampled points. Since by Rademacher’s theorem and countable additivity of probabilities this holds for every sampled point with probability one, we do not include a termination condition here.

---

**Algorithm GS** (Gradient Sampling)

---

**Require:** initial point  $x^0$  at which  $f$  is differentiable, initial sampling radius  $\epsilon_0 \in (0, \infty)$ , initial stationarity target  $\nu_0 \in [0, \infty)$ , sample size  $m \geq n + 1$ , line search parameters  $(\beta, \gamma) \in (0, 1) \times (0, 1)$ , termination tolerances  $(\epsilon_{\text{opt}}, \nu_{\text{opt}}) \in [0, \infty) \times [0, \infty)$ , and reduction factors  $(\theta_\epsilon, \theta_\nu) \in (0, 1] \times (0, 1]$

- 1: **for**  $k \in \mathbb{N}$  **do**
- 2:     independently sample  $\{x^{k,1}, \dots, x^{k,m}\}$  uniformly from  $\mathbb{B}(x^k, \epsilon_k)$
- 3:     compute  $g^k$  as the solution of  $\min_{g \in \mathcal{G}^k} \frac{1}{2} \|g\|_2^2$ , where  $\mathcal{G}^k := \text{conv}\{\nabla f(x^k), \nabla f(x^{k,1}), \dots, \nabla f(x^{k,m})\}$
- 4:     **if**  $\|g^k\|_2 \leq \nu_{\text{opt}}$  and  $\epsilon_k \leq \epsilon_{\text{opt}}$  **then** terminate
- 5:     **if**  $\|g^k\|_2 \leq \nu_k$
- 6:         **then** set  $\nu_{k+1} \leftarrow \theta_\nu \nu_k$ ,  $\epsilon_{k+1} \leftarrow \theta_\epsilon \epsilon_k$ , and  $t_k \leftarrow 0$
- 7:         **else** set  $\nu_{k+1} \leftarrow \nu_k$ ,  $\epsilon_{k+1} \leftarrow \epsilon_k$ , and
 
$$t_k \leftarrow \max \{t \in \{1, \gamma, \gamma^2, \dots\} : f(x^k - t g^k) < f(x^k) - \beta t \|g^k\|_2^2\} \quad (2.1)$$
- 8:     **if**  $f$  is differentiable at  $x^k - t_k g^k$
- 9:         **then** set  $x^{k+1} \leftarrow x^k - t_k g^k$
- 10:        **else** set  $x^{k+1}$  randomly as any point where  $f$  is differentiable such that
 
$$f(x^{k+1}) < f(x^k) - \beta t_k \|g^k\|_2^2 \quad \text{and} \quad \|x^k - t_k g^k - x^{k+1}\|_2 \leq \min\{t_k, \epsilon_k\} \|g^k\|_2 \quad (2.2)$$

11: **end for**

---

5. Finally, Steps 8–10 in Algorithm GS do require explicit checks that ensure that  $f$  is differentiable at  $x^{k+1}$ , but unlike in the variants in [BLO05] and [Kiw07], it is not required that  $f$  be *continuously* differentiable at  $x^{k+1}$ . This differentiability requirement is included since it is not the case that  $f$  is differentiable at  $x^k - t_k g^k$  with probability one, as is shown via an example in [HSS16], discussed further in §4.2. For a precise procedure for implementing Step 10, see [Kiw07].

The computation in Step 3 requires solving a strongly convex quadratic optimization problem (QP) to compute the minimum-norm element of the convex hull of the current and sampled gradients, or, equivalently, to compute the 2-norm projection of the origin onto this convex hull. It is essentially the same operation required in every iteration of a bundle method. To see this, observe that solving the QP in Step 3 can be expressed, with

$$G_k := [\nabla f(x^k) \quad \nabla f(x^{k,1}) \quad \dots \quad \nabla f(x^{k,m})],$$

as computing  $(z_k, d^k, \lambda^k) \in \mathbb{R} \times \mathbb{R}^n \times \mathbb{R}^{m+1}$  as the primal-dual solution of the QPs

$$\left\{ \begin{array}{l} \min_{(z,d) \in \mathbb{R} \times \mathbb{R}^n} z + \frac{1}{2} \|d\|_2^2 \\ \text{s.t. } G_k^T d \leq z \mathbf{1} \end{array} \right\} \quad \text{and} \quad \left\{ \begin{array}{l} \min_{\lambda \in \mathbb{R}^{m+1}} \frac{1}{2} \|G_k \lambda\|_2^2 \\ \text{s.t. } \mathbf{1}^T \lambda = 1, \lambda \geq 0 \end{array} \right\}. \quad (2.3)$$

The latter problem, yielding  $G_k \lambda^k = g^k$ , can easily be seen to be equivalent to solving the subproblem  $\min_{g \in \mathcal{G}^k} \frac{1}{2} \|g\|_2^2$  stated in Step 3, whereas the former problem, yielding  $d^k = -g^k$ , can be seen to have the same form as the subproblems arising in bundle methods.

Normally, the initial stationarity target  $\nu_0$  is chosen to be positive and the reduction factors  $\theta_\nu$  and  $\theta_\epsilon$  are chosen to be less than one so that the stationarity target and sampling radius are reduced every time the condition  $\|g^k\| \leq \nu_k$  is satisfied. However, it is also interesting to consider the variant with  $\nu_0 = 0$  and  $\theta_\epsilon = 1$ , forcing the algorithm to run forever with  $\epsilon$  fixed unless it terminates with  $g^k = 0$  for some  $k \in \mathbb{N}$ . We consider both of these variants in the global convergence theory given in the next section.

### 3 Convergence Theory for Algorithm GS

We begin with some definitions. In the locally Lipschitz case, the Clarke subdifferential of  $f$  at  $x \in \mathbb{R}^n$  is defined by the convex hull of the limits of gradients of  $f$  on sequences converging to  $x$  [Cla75, Def. 1.1], i.e.,

$$\bar{\partial}f(x) = \text{conv} \left\{ \lim_{j \rightarrow \infty} \nabla f(y^j) : \{y^j\} \rightarrow x \text{ where } f \text{ is differentiable at } y^j \text{ for all } j \in \mathbb{N} \right\}. \quad (3.1)$$

A point  $x$  is Clarke stationary for  $f$  if  $0 \in \bar{\partial}f(x)$ . A more “robust” sense of stationarity with respect to  $f$  can be defined by considering the limits of gradients corresponding to limiting points *near*  $x$ ; in particular, given a radius  $\epsilon \in [0, \infty)$ , the Clarke  $\epsilon$ -subdifferential [Gol77] is given by<sup>2</sup>

$$\bar{\partial}_\epsilon f(x) = \text{conv} \{ \bar{\partial}f(\mathbb{B}(x, \epsilon)) \}. \quad (3.2)$$

A point  $x$  is Clarke  $\epsilon$ -stationary for  $f$  if  $0 \in \bar{\partial}_\epsilon f(x)$ . For all practical purposes, one cannot generally evaluate  $\bar{\partial}f$  (or  $\bar{\partial}_\epsilon f$ ) at (or near) any point where  $f$  is not differentiable. That said, Algorithm GS is based on the idea that one can approximate the minimum norm element in  $\bar{\partial}_\epsilon f(x^k)$  through random sampling of gradients in the ball  $\mathbb{B}(x^k, \epsilon)$ . To a large extent this idea is motivated by [BLO02a] which investigates how well the entire Clarke subdifferential  $\bar{\partial}f(x)$  can be approximated through random sampling. However, the results in [BLO02a] cannot be directly exploited in the analysis of the GS algorithm because the gradients are sampled only at a finite number of points near any given iterate.

#### 3.1 Global Convergence Guarantees

A critical aspect of theoretical convergence guarantees for Algorithm GS concerns the set of points where  $f$  is *continuously* differentiable, which we denote by  $D$ . Consideration of  $D$  played a crucial role in the analysis in both [BLO05] and [Kiw07], but there were some oversights concerning both the requirements of the algorithm with respect to  $D$  and the assumptions on  $D$ . Regarding the requirements of the algorithm with respect to  $D$ , there is actually no need, from a theoretical point of view, for either the iterates  $\{x^k\}$  or the randomly generated sampled points  $\{x^{k,j}\}$  to lie in  $D$ ; all that is needed is that  $f$  is differentiable at these points. Most implementations of GS algorithms do not attempt to check any form of differentiability in any case, but if one were to attempt to implement such a check, it is certainly more tractable to check for differentiability than continuous differentiability. Regarding the assumptions on  $D$ , in the theorems that we state below, we assume that  $D$  is an open set with full measure in  $\mathbb{R}^n$ . In contrast, the relevant assumption stated in [BLO05, Kiw07] is weaker, namely, that  $D$  is an open dense subset of  $\mathbb{R}^n$ . However, the proofs of convergence actually require the full measure assumption on  $D$  that we include below.<sup>3</sup>

There are three types of global convergence guarantees of interest for Algorithm GS: one when the input parameters ensure that  $\{\epsilon_k\} \searrow 0$ , one when  $\epsilon_k$  is repeatedly reduced but a positive stopping tolerance prevents it from converging to zero, and one when  $\epsilon_k = \epsilon > 0$  for all  $k$ . These lead to different properties for the iterate sequence. The first theorem below relates to cases when the stationarity tolerance and sampling radius tolerance are both set to zero so that the algorithm can never terminate.

**Theorem 3.1.** *Suppose that  $f$  is locally Lipschitz in  $\mathbb{R}^n$  and continuously differentiable on an open set  $D$  with full measure in  $\mathbb{R}^n$ . Suppose further that Algorithm GS is run with  $\nu_0 > 0$ ,  $\nu_{\text{opt}} = \epsilon_{\text{opt}} = 0$ , and strict reduction factors  $\theta_\nu < 1$  and  $\theta_\epsilon < 1$ . Then, with probability one, Algorithm GS is well defined in the sense that the gradients in Step 3 exist in every iteration, the algorithm does not terminate, and either*

- (i)  $f(x^k) \searrow -\infty$ , or
- (ii)  $\nu_k \searrow 0$ ,  $\epsilon_k \searrow 0$ , and every cluster point of  $\{x^k\}$  is Clarke stationary for  $f$ .

Theorem 3.1 is essentially the same as [Kiw07, Thm 3.3] (with the modifications given in [Kiw07, §4.1] for nonnormalized directions), except for two aspects:

<sup>2</sup>The definition in [BLO05] includes a closure operation but this is unnecessary.

<sup>3</sup>This oversight went unnoticed for 12 years until J. Portegies and T. Mitchell brought it to our attention recently.

1. The proof given in [Kiw07] implicitly assumes that  $D$  is an open set with full measure, as does the proof of [BLO05, Thm 3.4] on which Kiwiel’s proof is based, although the relevant assumption on  $D$  in both papers is the weaker condition that  $D$  is an open dense set. Details are given in Appendix A.
2. In the algorithms analyzed in [Kiw07] and [BLO05], the iterates  $\{x^k\}$  and the randomly sampled points  $\{x^{k,j}\}$  were enforced to lie in the set  $D$  where  $f$  is continuously differentiable. We show in Appendix B that the theorem still holds when this requirement is relaxed to ensure only that  $f$  is differentiable at these points.

As Kiwiel argues, Theorem 3.1 is essentially the best result that could be expected. Furthermore, as pointed out in [Kiw07, Remark 3.7(ii)], it leads immediately to the following corollary.

**Corollary 3.1.** *Suppose that  $f$  is locally Lipschitz in  $\mathbb{R}^n$  and continuously differentiable on an open set  $D$  with full measure in  $\mathbb{R}^n$ . Suppose further that Algorithm GS is run with  $\nu_0 > \nu_{\text{opt}} > 0$ ,  $\epsilon_0 > \epsilon_{\text{opt}} > 0$ , and strict reduction factors  $\theta_\nu < 1$  and  $\theta_\epsilon < 1$ . Then, with probability one, Algorithm GS is well defined in the sense that the gradients in Step 3 exist at every iteration, and either*

- (i)  $f(x^k) \searrow -\infty$ , or
- (ii) Algorithm GS terminates by the stopping criteria in Step 4.

The final result that we state concerns the case when the sampling radius is fixed. A proof of this result is essentially given by that of [Kiw07, Thm 3.5], again taking into account the comments in Appendices A and B.

**Theorem 3.2.** *Suppose that  $f$  is locally Lipschitz in  $\mathbb{R}^n$  and continuously differentiable on an open set  $D$  with full measure in  $\mathbb{R}^n$ . Suppose further that Algorithm GS is run with  $\nu_0 = \nu_{\text{opt}} = 0$ ,  $\epsilon_0 = \epsilon_{\text{opt}} = \epsilon > 0$ , and  $\theta_\epsilon = 1$ . Then, with probability one, Algorithm GS is well defined in the sense that the gradients in Step 3 exist at every iteration, and one of the following occurs:*

- (a)  $f(x^k) \searrow -\infty$ , or
- (b) Algorithm GS terminates for some  $k \in \mathbb{N}$  with  $g^k = 0$ , or
- (c) there exists  $\mathcal{K} \subseteq \mathbb{N}$  with  $\{g^k\}_{k \in \mathcal{K}} \rightarrow 0$  and every cluster point of  $\{x^k\}_{k \in \mathcal{K}}$  is Clarke  $\epsilon$ -stationary for  $f$ .

Of the five open questions regarding the convergence analysis for gradient sampling raised in [BLO05], three were answered explicitly by Kiwiel in [Kiw07]. Another open question was: “Under what conditions can one guarantee that the GS algorithm terminates finitely?” This was posed in the context of a fixed sampling radius and therefore asks how one might know whether outcome (b) or (c) occurs in Theorem 3.2, assuming  $f$  is bounded below. This remains open, but Kiwiel’s introduction of the termination tolerances in the GS algorithm statement led to Corollary 3.1 which guarantees that when the sampling radius is reduced dynamically and the tolerances are nonzero, Algorithm GS must terminate if  $f$  is bounded below. The only other open question concerns extending the convergence analysis to the non-Lipschitz case.

Overall, Algorithm GS has a very satisfactory convergence theory in the locally Lipschitz case. Its main weakness is its per-iteration cost, most notably due to the need to compute  $m \geq n + 1$  gradients in every iteration and solve a corresponding QP. However, enhancements to the algorithm have been proposed that can vastly reduce this per-iteration cost while maintaining these guarantees. We discuss these and other enhancements in §4.

### 3.2 A Local Linear Convergence Result

Given the relationship between gradient sampling and a traditional steepest descent approach, one might ask if there are scenarios in which Algorithm GS can attain a linear rate of local convergence. The study in [HSS17] answers this in the affirmative, at least in a certain probabilistic sense. If (i) the set of sampled points is *good* in a certain sense described in [HSS17], (ii) the objective function  $f$  belongs to a class of

functions defined by the maximum of a finite number of smooth functions (“finite-max” functions), and (iii) the input parameters are set appropriately, then Algorithm **GS** will produce a step yielding a reduction in  $f$  that is significant. This analysis involves  $\mathcal{VU}$ -decomposition ideas [LOS00, Lew02, MS05], where in particular it is shown that the reduction in  $f$  is comparable to that achieved by a steepest descent method restricted to the smooth  $\mathcal{U}$ -space of  $f$ . This means that a linear rate of local convergence can be attained over any infinite subsequence of iterations in which the sets of sampled points are *good*.

### 3.3 The Non-Lipschitz Case

In the non-locally Lipschitz case, the Clarke subdifferential  $\bar{\partial}f$  is defined in [BLO02a, p. 573]; unlike in the Lipschitz case, this set may be unbounded, presenting obvious difficulties for approximating it through random sampling of gradients. In fact, more than half of [BLO02a] is devoted to investigating this issue, relying heavily on modern variational analysis as expounded in [Roc98]. Some positive results were obtained, specifically in the case that  $f$  is “directionally Lipschitz” at  $\bar{x}$ , which means that the “horizon cone” [BLO02a, p. 572] of  $f$  at  $\bar{x}$  is pointed, that is, it does not contain a line. For example, this excludes the function on  $\mathbb{R}$  defined by  $f(x) = |x|^{1/2}$  at  $\bar{x} = 0$ , but it applies to the case  $f(x) = (\max(0, x))^{1/2}$  even at  $\bar{x} = 0$ . The discussion of the directionally Lipschitz case culminates with Corollary 6.1, which establishes that the Clarke subdifferential can indeed be approximated by convex hulls of gradients. On the more negative side, Example 7.2 shows that this approximation can fail badly in the general Lipschitz case. Motivated by these results, Burke and Lin have recently extended the GS convergence theory to the directionally Lipschitz case [BL18, Lin09]. However, it would seem difficult to extend these results to the more general non-Lipschitz case.

Suppose  $f : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}$  is defined by

$$f(X) = \max\{\operatorname{Re} \lambda : \det(\lambda I - X) = 0\},$$

the spectral abscissa (maximum of the real parts of the eigenvalues) of  $X$ . Assume that  $\bar{X}$  has the property that its only eigenvalues whose real parts coincide with  $f$  make up a zero eigenvalue with multiplicity  $q$  associated with a single Jordan block (the generic case). In this case the results in [BO01] tell us that the horizon cone of  $f$  is pointed at  $\bar{X}$  if and only if the multiplicity  $q \leq 2$ ; on the other hand  $f$  is locally Lipschitz at  $X$  if and only if  $q = 1$ . In the light of the previous paragraph, one might expect much greater practical success in applying GS to minimize the spectral abscissa of a parameterized matrix if the optimal multiplicities are limited to 1 or 2. However, this seems not to be the case. The results reported in [BLO02b] for unconstrained spectral abscissa minimization, as well as results for applying the algorithm of [CO12] (see §5.2 below) for constrained nonsmooth, nonconvex optimization to problems with spectral radius objective and constraints, as reported in [CMO17, Sec. 4.2 and Appendix A.1], do not show any marked deterioration as the optimal multiplicities increase from 2 or 3, although certainly the problems are much more challenging for larger multiplicities. We view understanding the rather remarkably good behavior of the GS algorithm on such examples as a potentially rewarding, though certainly challenging, line of investigation.

## 4 Enhancements

As explained above, the statement of Algorithm **GS** differs in several ways from the algorithms stated in [BLO02b], [BLO05], and [Kiw07]. Other variants of the strategy have also been proposed in recent years, in some cases to pose new solutions to theoretical obstacles (such as the need, in theory, to check for differentiability of  $f$  at each new iterate), and in others to enhance the practical performance of the approach. In this section, we discuss a few of these enhancements.

### 4.1 Restricting the Line Search to within a Trust Region

Since the gradient information about  $f$  is obtained only within the ball  $\mathbb{B}(x^k, \epsilon_k)$  for all  $k \in \mathbb{N}$ , and since one might expect that smaller steps should be made when the sampling radius is small, an argument can be

made that the algorithm might benefit by restricting the line search to within the ball  $\mathbb{B}(x^k, \epsilon_k)$  for all  $k \in \mathbb{N}$ . In [Kiw07, §4.2], such a variant is proposed where in place of  $-g^k$  the search direction is set as  $-\epsilon_k g^k / \|g^k\|_2$ . With minor corresponding changes to conditions (2.1) and (2.2), all of the theoretical convergence guarantees of the algorithm are maintained. Such a variant with the trust region radius defined as a positive multiple of the sampling radius  $\epsilon_k$  for all  $k \in \mathbb{N}$  would have similar properties. This variant might perform well in practice, especially in situations when otherwise setting the search direction as  $-g^k$  would lead to significant effort being spent in the line search.

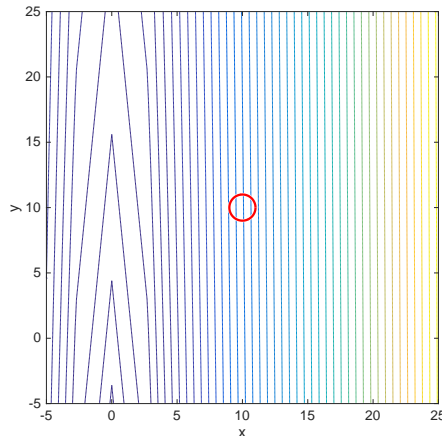
## 4.2 Avoiding the Differentiability Check

The largest distraction from the fundamentally simple nature of Algorithm GS is the procedure for choosing a perturbed subsequent iterate if  $f$  is not differentiable at  $x^k - t_k g^k$ ; see Steps 8–10. This procedure is necessary for the algorithm to be well defined since, to ensure that  $-g^k$  is a descent direction for all  $k \in \mathbb{N}$ , the algorithm relies on the existence of and ability to compute  $-\nabla f(x^k)$  for all  $k \in \mathbb{N}$ . One might hope that this procedure, while necessary for theoretical convergence guarantees, could be ignored in practice. However, due to the deterministic nature of the line search, situations exist in which landing on a point of nondifferentiability of  $f$  occurs with positive probability.

For example, consider the function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  given by

$$f(w, z) = \max\{0.5w^2 + 0.1z, w + 0.1z + 1, -w + 0.1z + 1, -0.05z - 50\};$$

see Figure 2. As shown by Helou, Santos, and Simões in [HSS16], if Algorithm GS is initialized at  $x^0 = (w_0, z_0)$  chosen anywhere in the unit ball centered at  $(10, 10)$ , then there is a positive probability that the function  $f$  will not be differentiable at  $x^0 - t_0 g^0$ . This can be explained as follows. At any point in the unit ball centered at  $(10, 10)$ , the function  $f$  is continuously differentiable and  $\nabla f(x^0) = [w_0; 0.1]^T$ . Moreover, there is a positive probability that the sampled points obtained at this first iteration will yield  $g^0 = \nabla f(x^0)$ . Therefore, given a reasonable value for the parameter  $\beta$  that appears in (2.1) (e.g.,  $\beta = 10^{-4}$ ), the sufficient decrease of the function value is attained with  $t_0 = 1$ . This guarantees that the function  $f$  will not be differentiable at the next iterate, since the first coordinate of  $x^1 = (w_1, z_1)$  will be zero.



**Figure 2** – Contours of a function illustrating the necessity of the differentiability check in Algorithm GS. Initialized uniformly within the illustrated ball, there is a positive probability that  $x^0 - t_0 g^0 \notin D$ .

The authors of [HSS16] propose two strategies to avoid the issues highlighted by this example. The first is that, rather than perturb the iterate after the line search, one could perturb the search direction before the line search. It is shown that if the random perturbation of the search direction is sufficiently small such that, for one thing, the resulting direction is still one of descent for  $f$ , then  $f$  will be differentiable at all



iterates with probability one. Their second proposed strategy involves the use of a nonmonotone line search. In particular, it is shown that if a strictly positive value  $\Delta_k$  is added on the right-hand side of the sufficient decrease condition in (2.1) such that  $\{\Delta_k\}$  is summable, then one can remove  $\nabla f(x^k)$  from the set  $\mathcal{G}_k$  for all  $k \in \mathbb{N}$  and maintain convergence guarantees even when  $f$  is *not* differentiable at  $\{x^k\}$ . This can be shown by noticing that, due to the positive term  $\Delta_k$ , the line search procedure continues to be well defined even if  $-g^k$  is not a descent direction for  $f$  at  $x^k$  (which may happen since  $\nabla f(x_k)$  is no longer involved in the computation of  $g^k$ ). However, the summability of  $\{\Delta_k\}$  implies that the possible increases in  $f$  will be finite and, when the sampled points are good enough to describe the local behavior of  $f$  around the current iterate, the function value will necessarily decrease if the method has not reached (approximate) stationarity. Overall, the sufficient reductions in  $f$  achieved in certain iterations will ultimately exceed any increases.

Another proposal for avoiding the need to have  $f$  differentiable at  $x^k$  for all  $k \in \mathbb{N}$  is given in [Kiw07, §4.3], wherein a technique is proposed for using a limited line search, potentially causing the algorithm to take null steps in some iterations. In fact, this idea of using a limited line search can also be used to avoid the need to sample a new set of  $m \geq n + 1$  gradients in each iteration, as we discuss next.

### 4.3 Adaptive Sampling

As previously mentioned, the main weakness of Algorithm GS is the cost of computing  $m \geq n + 1$  gradients in every iteration and solving a corresponding QP. Loosely speaking, this lower bound on the number of gradients is required in the analysis of the algorithm so that one can use Carathéodory’s theorem to argue that, with at least  $n + 1$  gradients, there is a sufficiently good chance that the combination vector  $-g^k$  represents a sufficiently good direction from  $x^k$ . However, in many situations in practice, the sampling of such a large number of gradients in each iteration can lead to a significant amount of wasted computational effort. One can instead sample *adaptively*, attempting to search along directions computed using fewer gradients and proceeding as long as a sufficient reduction is attained.

In [CQ13], Curtis and Que show how such an adaptive sampling strategy can be employed so that the convergence guarantees of Algorithm GS are maintained while only a constant number (independent of  $n$ ) of gradients need to be sampled in each iteration. A key aspect that allows one to maintain these guarantees is the employment of a limited line search, as first proposed in [Kiw07], potentially leading to a null step when fewer than  $n + 1$  gradients are currently in hand and when the line search is not successful after a prescribed finite number of function evaluations. See also [CQ15] for further development of these ideas, where it is shown that one might not need to sample *any* gradients as long as a sufficient reduction is attained.

The work in [CQ13] also introduces the idea that, when adaptive sampling is employed, the algorithm can exploit a practical feature commonly used in bundle methods. This idea relates to warm-starting the algorithm for solving the QP subproblems. In particular, suppose that one has solved the primal-dual pair of QPs in (2.3) for some  $m \in \mathbb{N}$  to obtain  $(z_k, d^k, \lambda^k)$ , yielding  $g^k = -d^k$ . If one were to subsequently aim to solve the pair of QPs corresponding to the augmented matrix of gradients

$$\overline{G}_k = [G_k \quad \nabla f(x^{k,m+1}) \quad \dots \quad \nabla f(x^{k,m+p})],$$

then one obtains a viable feasible starting point for the latter QP in (2.3) by augmenting the vector  $\lambda^k$  with  $p$  zeros. This can be exploited, e.g., in an active-set method for solving this QP; see [Kiw86].

As a further practical enhancement, the work in [CQ13, CQ15] also proposes the natural idea that, after moving to a new iterate  $x^k$ , gradients computed in previous iterations can be “re-used” if they correspond to points that lie within  $\mathbb{B}(x^k, \epsilon_k)$ . This may further reduce the number of sampled points needed in practice.

### 4.4 Second-Order-Type Variants

The solution vector  $d^k$  of the QP in (2.3) can be viewed as the minimizer of the model of  $f$  at  $x^k$  given by

$$q_k(d) = f(x^k) + \max_{g \in \mathcal{G}_k} g^T d + \frac{1}{2} d^T H_k d$$

with  $H_k = I$ , the identity matrix. As in other second-order-type methods for nonlinear optimization, one might also consider algorithm variants where  $H_k$  is set to some other symmetric positive definite matrix. Ideas of this type have been explored in the literature. For example, in [CQ13], two techniques are proposed: one in which  $H_k$  is set using a quasi-Newton updating strategy and one in which the matrix is set in an attempt to ensure that the model  $q_k$  represents an upper bounding model for  $f$ . The idea of employing a quasi-Newton approach, inspired by the success of quasi-Newton methods in practice for nonsmooth optimization (see [LO13]), has also been explored further in [CQ15, CRZ17].

Another approach, motivated by the encouraging results obtained when employing spectral gradient methods to solve smooth [Fle05, BMR14] and nonsmooth [CLR07] optimization problems, has been to employ a Barzilai-Borwein (BB) strategy for computing initial stepsizes in a GS approach; see [LACR17] and the background in [BB88, Ray93, Ray97]. Using a BB strategy can be viewed as choosing  $H_k = \alpha_k I$  for all  $k \in \mathbb{N}$  where the scalar  $\alpha_k$  is set according to iterate and gradient displacements in the latest iteration.

In all of these second-order-type approaches, one is able to maintain convergence guarantees of the algorithm as long as the procedure for setting the matrix  $H_k$  is safeguarded during the optimization process. For example, one way to maintain guarantees is to restrict each  $H_k$  to the set of symmetric matrices whose eigenvalues are contained within a fixed positive interval. One might also attempt to exploit the self-correcting properties of BFGS updating; see [CRZ17].

## 5 Extensions

In this section, we discuss extensions to the GS methodology to solve classes of problems beyond unconstrained optimization on  $\mathbb{R}^n$ .

### 5.1 Riemannian GS for Optimization on Manifolds

Hosseini and Uschmajew in [HU17] have extended the GS methodology for minimizing a locally Lipschitz  $f$  over a set  $\mathcal{M}$ , where  $\mathcal{M}$  is a complete Riemannian manifold of dimension  $n$ . The main idea of this extension is to employ the convex hull of gradients from tangent spaces at randomly sampled points *transported* to the tangent space of the current iterate. In this manner, the algorithm can be characterized as a generalization of the Riemannian steepest descent method just as GS is a generalization of traditional steepest descent. Assuming that the vector transport satisfies certain assumptions, including a *locking condition*, the algorithm attains convergence guarantees on par with those for Algorithm GS.

### 5.2 SQP-GS for Constrained Optimization

Curtis and Overton in [CO12] proposed a combination sequential quadratic programming (SQP) and gradient sampling method for solving constrained optimization problems involving potentially nonsmooth constraint functions, i.e.,

$$\min_{x \in \mathbb{R}^n} f(x) \quad \text{s.t.} \quad c(x) \leq 0, \quad (5.1)$$

where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  and  $c : \mathbb{R}^n \rightarrow \mathbb{R}^{n_c}$  are locally Lipschitz. A key aspect of the proposed SQP-GS approach is that sampled points for the objective and each individual constraint function are generated independently. With this important feature, it is shown that the algorithm, which follows a penalty-SQP strategy (e.g., see [Fle87]), attains convergence guarantees for minimizing an exact penalty function that are similar to those in §3.1. Moreover, with the algorithm’s penalty parameter updating strategy, it is shown that either the penalty function is driven to  $-\infty$ , the penalty parameter settles at a finite value and any limit point will be feasible for the constraints and stationary for the penalty function, or the penalty parameter will be driven to zero and any limit point of the algorithm will be stationary for a constraint violation measure. As for other exact penalty function methods for nonlinear optimization, one can translate between guarantees for minimizing the exact penalty function and solving the constrained problem (5.1); in particular, if the problem is *calm* [Cla83, Roc82], then at any local minimizer  $x_*$  of (5.1) there exists a threshold for the penalty parameter beyond which  $x_*$  will be a local minimizer of the penalty function.

Tang, Liu, Jian, and Li have also proposed in [TLJL14] a *feasible* variant of the SQP-GS method in which the iterates are forced to remain feasible for the constraints and the objective function is monotonically decreasing throughout the optimization process. This opens the door to employing a two-phase approach common for solving some optimization problems, where phase 1 is responsible for attaining a feasible point and phase 2 seeks optimality while maintaining feasibility.

### 5.3 Derivative-Free Optimization

Given its simple nature, gradient sampling has proved to be an attractive basis for the design of new algorithms even when gradient information cannot be computed explicitly. Indeed, there have been a few variants of *derivative-free* algorithms that have been inspired by gradient sampling.

The first algorithm for derivative-free optimization inspired by gradient sampling was proposed by Kiwiel in [Kiw10]. In short, in place of the gradients appearing in Algorithm GS, this approach employs Gupal’s estimates of gradients of the Steklov averages of  $f$ . In this manner, function values only—specifically,  $\mathcal{O}(mn)$  per iteration—are required for convergence guarantees. A less expensive *incremental* version is also proposed.

Another derivative-free variant of GS, proposed by Hare and Nutini in [HN13], is specifically designed for minimizing finite-max functions. This approach exploits knowledge about which of these functions are *almost active*—in terms of having value close to the objective function—at a particular point. In so doing, rather than attempt to approximate gradients at nearby points, as in done in [Kiw10], this approach only attempts to approximate gradients of almost active functions. The convergence guarantees proved for the algorithm are similar to those for GS methods, though the practical performance is improved by the algorithm’s tailored gradient approximation strategy.

Finally, we mention the *manifold sampling* algorithm, proposed by Larson, Menickelly, and Wild in [LMW16], for solving nonconvex problems where the objective function is the  $\ell_1$ -norm of a smooth vector function  $F : \mathbb{R}^n \rightarrow \mathbb{R}^r$ . While this approach does not employ a straightforward GS methodology in that it does not randomly sample points, it does employ a GS-type approach in the way that the gradient of a model of the objective function is constructed by solving a QP of the type in (2.3). Random sampling can be avoided in this construction since the algorithm can exploit knowledge of the signs of the elements of  $F(x)$  at any  $x \in \mathbb{R}^n$  along with knowledge of  $\|\partial\| \cdot \|\cdot\|_1$ .

## 6 Applications

We mentioned in the introduction that the original gradient sampling paper [BLO02b] reported results for spectral abscissa optimization problems that had not been solved previously. The second gradient sampling paper [BLO05] reported results for many more applications that again had not been solved previously: these included Chebyshev approximation by exponential sums, eigenvalue product minimization for symmetric matrices, spectral and pseudospectral abscissa minimization, maximization of the “distance to instability”, and fixed-order controller design by static output feedback.

Subsequently, the GS algorithm played a key role in the HANSO (Hybrid Algorithm for Non-Smooth Optimization)<sup>4</sup> and HIFOO (H-Infinity Fixed-Order Optimization)<sup>5</sup> toolboxes. The former is a stand-alone code for unconstrained nonsmooth optimization while the latter is a more specialized code used for the design of low-order controllers for linear dynamical systems with input and output, computing fixed-order controllers by optimizing stability measures that are generally nonsmooth at local minimizers [BHLO06]. HIFOO calls HANSO to carry out the optimization. The use of “hybrid” in the expansion of the HANSO acronym indicated that, from its inception, HANSO combined the use of both a quasi-Newton algorithm (BFGS) and Gradient Sampling. The quasi-Newton method was used in an initial phase which, rather surprisingly, typically worked very effectively even in the presence of nonsmoothness, very often providing a fast way to approximate a local minimizer. This was followed by a GS phase to refine the approximation, typically verifying a loose measure of local optimality. The HIFOO toolbox has been used successfully in a wide

<sup>4</sup>[www.cs.nyu.edu/overton/software/hanso/](http://www.cs.nyu.edu/overton/software/hanso/)

<sup>5</sup>[www.cs.nyu.edu/overton/software/hifoo/](http://www.cs.nyu.edu/overton/software/hifoo/)

variety of applications, including synchronization of heterogeneous multi-agent systems and networks, design of motorized gimbals that stabilize an angular motion of an optical payload around an axis, flight control via static output feedback, robust observer-based fault detection and isolation, influence of tire damping on control of quarter-car suspensions, flexible aircraft lateral flight dynamic control, optimal control of aircraft with a blended wing body, vibration control of a fluid/plate system, controller design of a nose landing gear steering system, bilateral teleoperation for minimally invasive surgery, design of an aircraft controller for improved gust alleviation and passenger comfort, robust controller design for a proton exchange membrane fuel cell system, design of power systems controllers, and design of winding systems for elastic web materials — for a full list of references, see [CMO17].

The successful use of BFGS in HANSO and HIFOO led to papers on the use of quasi-Newton methods in the nonsmooth context, both for unconstrained [LO13] and constrained [CMO17] optimization. The latter paper introduced a new BFGS-SQP method for nonsmooth constrained optimization and compared it with the SQP-GS method discussed in § 5.2 on a suite of challenging static output feedback controller design problems, half of them non-Lipschitz (spectral radius minimization) and half of them locally Lipschitz (pseudospectral radius minimization). It was found that although the BFGS-SQP method was much faster than SQP-GS, nonetheless, if the latter method based on Gradient Sampling was allowed sufficient running time, it frequently found better approximate solutions than the former method based on BFGS in a well defined sense, evaluated using “relative minimization profiles”. Interestingly, this was particularly pronounced on the non-Lipschitz problems, despite the fact that the GS convergence theory does not extend to this domain. See §3.3 for further discussion of this issue.

Finally, we mention an interesting application of GS to robot path planning [TM16]. This work is based on the observation that shortest paths generated through gradient descent on a value function have a tendency to chatter and/or require an unreasonable number of steps to synthesize. The authors demonstrate that the GS algorithm can largely alleviate this problem. For systems subject to state uncertainty whose state estimate is tracked using a particle filter, they proposed the Gradient Sampling with Particle Filter (GSPF) algorithm, which uses the particles as the locations in which to sample the gradient. At each step, the GSPF efficiently finds a consensus direction suitable for all particles or identifies the type of stationary point on which it is stuck. If the stationary point is a minimum, the system has reached its goal (to within the limits of the state uncertainty) and the algorithm terminates; otherwise, the authors propose two approaches to find a suitable descent direction. They illustrated the effectiveness of the GSPF on several examples using well known robot simulation environments. This work was recently extended and modified in [EM18], where the practical effectiveness of both the GSPF algorithm and the new modification was demonstrated on a Segway Robotic Mobility Platform.

## 7 Conclusion and Future Directions

Gradient sampling is a conceptually straightforward approximate steepest descent method. With a solid convergence theory, the method has blossomed into a powerful methodology for solving nonsmooth minimization problems. The theme of our treatment of GS in this work has been to emphasize the fact that, even though the basic algorithm has been enhanced and extended in various ways, the foundation of the approach is fundamentally simple in nature.

We have also corrected an oversight in the original GS theory (i.e., that the convergence results depend on assuming that the set of points over which the Lipschitz function  $f$  is continuously differentiable has full measure, although we do not have a counterexample to convergence of GS in the absence of this assumption). At the same time we have loosened the requirements of the algorithm (showing that  $f$  need only be differentiable at the iterates and sampled points). An open question that still remains is whether one can extend the GS theory to broader function classes, such as the case where  $f$  is assumed to be semi-algebraic but not necessarily locally Lipschitz or directionally Lipschitz.

Opportunities for extending GS theory for broader function classes may include connecting the algorithm to other randomized/stochastic optimization methods. For example, one might view the algorithm as a stochastic-gradient-like method applied to a smoothed objective. (A similar philosophy underlies the analysis

by Nesterov and Spokoiny in [NS17].) More precisely, given a locally Lipschitz objective  $f$ , consider a smoothing  $f_\epsilon$  whose value at any point  $x$  is given by the mean value of  $f$  over the ball  $\mathbb{B}(x, \epsilon)$ . The GS algorithm uses gradients of  $f$  at uniformly distributed random points in this ball. Notice that each such gradient can also be viewed as a stochastic gradient for the smoothing  $f_\epsilon$  in the sense that its expectation is the gradient of  $f_\epsilon$  at  $x$ . Thus, one might hope to prove convergence results for a GS algorithm (with predetermined stepsizes rather than line searches) that parallel convergence theory for stochastic gradient methods. Recent work by Davis, Drusvyatskiy, Kakade and Lee [DDKL18] gives convergence results for stochastic *subgradient* methods on a broad class of problems.

Another potentially interesting connection is with the work of Davis and Drusvyatskiy [DD18] on stochastic model-based optimization. Consider a GS variant that successively minimizes stochastic models of the objective function  $f$ , where we assume for simplicity that  $f$  is a globally Lipschitz convex function. In this variant, rather than moving along the direction  $-g^k$ , consider instead the construction of a cutting plane approximation of  $f$  from its affine minorants at the current iterate  $x^k$  and the sampled points  $\{x^{k,i}\}$ , augmented by the proximal term  $\beta_k \|x - x^k\|^2$ , where  $\{\beta_k\}$  is a predetermined sequence. Suppose that the next iterate is chosen as the minimizer of this model; for a given  $k$  and with  $\beta_k = 1$ , by equation (2.3), this scheme and GS produce similar descent directions as the sampling radius tends to zero. It follows from the results of [DD18] that the expected norm of the gradient of the Moreau envelope [DD18, p. 6] is reduced below  $\epsilon$  in  $\mathcal{O}(\epsilon^{-4})$  iterations. In fact, the assumptions on  $f$  in [DD18] are substantially weaker than convexity, and do not require any property of the set on which  $f$  is continuously differentiable.

Connecting the convergence theory for GS to stochastic methods as suggested in the previous two paragraphs could be enlightening. However, while stochastic methods are often designed for settings in which it is intractable to compute function values exactly—a feature reflected in the fact that the analyses for such methods are based on using predetermined stepsize sequences—the GS methodology has so far been motivated by problems for which functions and gradients are tractable to compute. In such settings, the line search in Algorithm GS is an ingredient that is crucial to its practical success.

## A On the Assumption of Full Measure of the Set $D$

Recall that  $D$  is defined to be the set of points on which the locally Lipschitz function  $f$  is continuously differentiable. There are two ways in which the analyses in [Kiw07, BLO05] actually depend on  $D$  having full measure:

1. The most obvious is that both papers require that the points sampled in each iteration should lie in  $D$ , and a statement is made in both papers that this occurs with probability one, but this is not the case if  $D$  is assumed only to be an open dense subset of  $\mathbb{R}^n$ . However, as already noted earlier and justified in Appendix B, this requirement can be relaxed, as in Algorithm GS given in §2, to require only that  $f$  be differentiable at the sampled points.
2. The set  $D$  must have full measure for Property 1, stated below, to hold. The proofs in [BLO05, Kiw07] depend critically on this property, which follows from [BLO02a, Eq. 1.2] (where it was stated without proof). For completeness we give a proof here, followed by an example that demonstrates the necessity of the full measure assumption.

**Property 1.** *Assume that  $D$  has full measure and let*

$$G_\epsilon(x) := \text{cl conv } \nabla f(\mathbb{B}(x, \epsilon) \cap D).$$

*For all  $\epsilon > 0$  and all  $x \in \mathbb{R}^n$ , one has  $\bar{\partial}f(x) \subseteq G_\epsilon(x)$ , where the Clarke subdifferential  $\bar{\partial}f$  is defined in (3.1).*

*Proof.* Let  $x \in \mathbb{R}^n$  and  $v \in \bar{\partial}f(x)$ . We have from [Cla83, Thm 2.5.1] that, for any set  $S$  with zero measure,

$$\bar{\partial}f(x) = \text{conv} \left\{ \lim_j \nabla f(y^j) : y^j \rightarrow x \text{ where } y^j \notin S \text{ and } f \text{ is differentiable at } y^j, \text{ for all } j \in \mathbb{N} \right\}.$$

In particular, since  $D$  has full measure and  $f$  is differentiable on  $D$ , it follows that

$$\bar{\partial}f(x) = \text{conv} \left\{ \lim_j \nabla f(y^j) : y^j \rightarrow x \text{ with } y^j \in D \text{ for all } j \in \mathbb{N} \right\}.$$

Considering this last relation and Carathéodory's theorem, it follows that  $v \in \text{conv} \{\xi^1, \dots, \xi^{n+1}\}$ , where, for all  $i \in \{1, \dots, n+1\}$ , one has  $\xi^i = \lim_j \nabla f(y^{j,i})$  for some sequence  $\{y^{j,i}\}_{j \in \mathbb{N}} \subset D$  converging to  $x$ . Hence, there must exist a sufficiently large  $j_i \in \mathbb{N}$  such that

$$y^{j,i} \in \mathbb{B}(x, \epsilon) \cap D \implies \nabla f(y^{j,i}) \in \nabla f(\mathbb{B}(x, \epsilon) \cap D) \subseteq \text{conv} \nabla f(\mathbb{B}(x, \epsilon) \cap D) \quad \text{for all } j \geq j_i.$$

Recalling that  $G_\epsilon(x)$  is the closure of  $\text{conv} \nabla f(\mathbb{B}(x, \epsilon) \cap D)$ , it follows that  $\xi^i \in G_\epsilon(x)$  for all  $i \in \{1, \dots, n+1\}$ . Moreover, since  $G_\epsilon(x)$  is convex, it follows that  $v \in G_\epsilon(x)$ . The result follows since  $x \in \mathbb{R}^n$  and  $v \in \bar{\partial}f(x)$  were arbitrarily chosen.  $\square$

With the assumption that  $D$  has full measure, Property 1 holds and hence the proofs of the results in [BLO05, Kiw07] are all valid. In particular, the proof of (ii) in [Kiw07, Lemma 3.2], which borrows from [BLO05, Lemma 3.2], depends on Property 1. See also the top of [BLO05, p. 762].

The following example shows that Property 1 might not hold if  $D$  is assumed only to be an open dense set, not necessarily of full measure.

**Example 1.** Let  $\delta \in (0, 1)$  and  $\{q_k\}_{k \in \mathbb{N}}$  be the enumeration of the rational numbers in  $(0, 1)$ . Define

$$D := \bigcup_{k=1}^{\infty} \mathcal{Q}_k, \text{ where } \mathcal{Q}_k := \left( q_k - \frac{\delta}{2^{k+1}}, q_k + \frac{\delta}{2^{k+1}} \right).$$

Clearly, its Lebesgue measure satisfies  $0 < \lambda(D) \leq \delta$ . Moreover, the set  $D$  is an open dense subset of  $[0, 1]$ . Now, let  $i_D : [0, 1] \rightarrow \mathbb{R}$  be the indicator function of the set  $D$ ,

$$i_D(x) = \begin{cases} 1, & \text{if } x \in D \\ 0, & \text{if } x \notin D \end{cases}.$$

Then, considering the Lebesgue integral, we define the function  $f : [0, 1] \rightarrow \mathbb{R}$ ,

$$f(x) = \int_{[0,x]} i_D d\lambda.$$

Let us prove that  $f$  is a Lipschitz continuous function on  $(0, 1)$ . To see this, note that given any  $a, b \in (0, 1)$  with  $b > a$ , it follows that

$$|f(b) - f(a)| = \left| \int_{[0,b]} i_D d\lambda - \int_{[0,a]} i_D d\lambda \right| = \left| \int_{(a,b)} i_D d\lambda \right| \leq \int_{(a,b)} 1 d\lambda = b - a,$$

which ensures that  $f$  is a Lipschitz continuous function on  $(0, 1)$ . Consequently, the Clarke subdifferential set of  $f$  at any point in  $(0, 1)$  is well defined. Moreover, we claim that, for all  $k \in \mathbb{N}$ ,  $f$  is continuously differentiable at any point  $q \in \mathcal{Q}_k$  and the following holds

$$f'(q) = i_D(q) = 1. \tag{1}$$

Indeed, given any  $q \in \mathcal{Q}_k$ , we have

$$f(q+t) - f(q) = \int_{[0,q+t]} i_D d\lambda - \int_{[0,q]} i_D d\lambda = \int_{(q,q+t)} i_D d\lambda, \text{ for } t > 0.$$

Since  $\mathcal{Q}_k$  is an open set, we can find  $\bar{t} > 0$  such that  $[q, q+t] \subset \mathcal{Q}_k \subset D$ , for all  $t \leq \bar{t}$ . Hence, given any  $t \in (0, \bar{t}]$ , it follows that

$$f(q+t) - f(q) = \int_{(q, q+t)} 1 d\lambda = t \implies \lim_{t \searrow 0} \frac{f(q+t) - f(q)}{t} = 1 = i_D(q).$$

The same reasoning can be used to see that the left derivative of  $f$  at  $q$  exists and it is equal to  $i_D(q)$ . Consequently, we have  $f'(q) = i_D(q) = 1$  for all  $q \in \mathcal{Q}_k$ , which yields that  $f$  is continuously differentiable on  $D$ .

By the Lebesgue differentiation theorem, we know that  $f'(x) = i_D(x)$  almost everywhere. Since the set  $[0, 1] \setminus D$  does not have measure zero, this means that there must exist  $z \in [0, 1] \setminus D$  such that  $f'(z) = i_D(z) = 0$ . Defining  $\epsilon := \min\{z, 1-z\}/2$ , we see, by (1), that the set

$$G_\epsilon(z) := \text{cl conv } \nabla f([z - \epsilon, z + \epsilon] \cap D)$$

is a singleton  $G_\epsilon(z) = \{1\}$ . However, since  $f'(z) = 0$ , it follows that  $0 \in \bar{\partial}f(z)$ , which implies  $\bar{\partial}f(z) \not\subseteq G_\epsilon(z)$ .

Note that it is stated on [BLO05, p. 754] and [Kiw07, p. 381] that the following holds: for all  $0 \leq \epsilon_1 < \epsilon_2$  and all  $x \in \mathbb{R}^n$ , one has  $\bar{\partial}_{\epsilon_1} f(x) \subseteq G_{\epsilon_2}(x)$ . Property 1 is a special case of this statement with  $\epsilon_1 = 0$ , and hence this statement too holds only under the full measure assumption.

Finally, it is worth mentioning that in practice, the full measure assumption on  $D$  usually holds. In particular, whenever a real-valued function is semi-algebraic (or, more generally, “tame”) — in other words, for all practical purposes virtually always — it is continuously differentiable on an open set of full measure. Hence, the original proofs hold in such contexts.

## B On Relaxing the Requirement that the Objective is Continuously Differentiable at the Iterates and Sampled Points

In this appendix, we summarize why it is not necessary that the iterates and sampled points of the algorithm lie in the set  $D$  in which  $f$  is continuously differentiable, and that rather it is sufficient to ensure that  $f$  is differentiable at these points, as in Algorithm GS. We do this by outlining how to modify the proofs in [Kiw07] to extend to this case.

1. That the gradients at the sampled points  $\{x^{k,j}\}$  exist follows with probability one from Rademacher’s theorem, while the existence of the gradients at the iterates  $\{x^k\}$  is ensured by the statement of Algorithm GS. Notice that the proof of part (ii) of [Kiw07, Theorem 3.3] still holds in our setting with the statement that the components of the sampled points are “sampled independently and uniformly from  $\mathbb{B}(x^k, \epsilon) \cap D$ ” replaced with “sampled independently and uniformly from  $\mathbb{B}(x^k, \epsilon)$ ”.
2. One needs to verify that  $f$  being differentiable at  $x^k$  is enough to ensure that the line search procedure presented in (2.1) terminates finitely. This is straightforward. Since  $\nabla f(x^k)$  exists, it follows that the directional derivative along any vector  $d \in \mathbb{R}^n \setminus \{0\}$  is given by  $f'(x^k; d) = \nabla f(x^k)^T d$ . Furthermore, since  $-\nabla f(x^k)^T g^k \leq -\|g^k\|_2^2$  (see [BLO05, p. 756]), it follows, for any  $\beta \in (0, 1)$ , that there exists  $\bar{t} > 0$  such that

$$f(x^k - tg^k) < f(x^k) - t\beta\|g^k\|_2^2 \quad \text{for any } t \in (0, \bar{t}).$$

This shows that the line search is well defined.

3. The only place where we actually need to modify the proof in [Kiw07] concerns item (ii) in Lemma 3.2, where it is stated that  $\nabla f(x^k) \in G_\epsilon(\bar{x})$  (for a particular point  $\bar{x}$ ) because  $x^k \in \mathbb{B}(\bar{x}, \epsilon/3) \cap D$ ; the latter is not true if  $x^k \notin D$ . However, using Property 1, we have

$$\nabla f(x^k) \in \bar{\partial}f(x^k) \subset G_{\epsilon/3}(x^k) \subset G_\epsilon(\bar{x}) \quad \text{when } x^k \in \mathbb{B}(\bar{x}, \epsilon/3),$$

so  $\nabla f(x^k) \in G_\epsilon(\bar{x})$  even when  $x^k \notin D$ .

Finally, although it was convenient in Appendix A to state Property 1 in terms of  $D$ , it actually holds if  $D$  is replaced by any full measure set on which  $f$  is differentiable. Nonetheless, it is important to note that the proofs of the results in [BLO05, Kiw07] do require that  $f$  be *continuously* differentiable on  $D$ . This assumption is used in the proof of (i) in [Kiw07, Lemma 3.2].

## References

- [AO17] A. Asl and M. L. Overton. Analysis of the Gradient Method with an Armijo-Wolfe Line Search on a Class of Nonsmooth Convex Functions, November 2017. arXiv:1711.08517.
- [BB88] J. Barzilai and J. M. Borwein. Two-Point Step Size Gradient Methods. *IMA Journal of Numerical Analysis*, 8(1):141–148, 1988.
- [BHLO06] J. V. Burke, D. Henrion, A. S. Lewis, and M. L. Overton. HIFOO - a MATLAB package for fixed-order controller design and  $H_\infty$  optimization. In *Fifth IFAC Symposium on Robust Control Design, Toulouse*, 2006.
- [BL18] J. V. Burke and Q. Lin. The Gradient Sampling Algorithm for Directionally Lipschitzian Functions, 2018. To appear.
- [BLO02a] J. V. Burke, A. S. Lewis, and M. L. Overton. Approximating Subdifferentials by Random Sampling of Gradients. *Math. Oper. Res.*, 27(3):567–584, 2002.
- [BLO02b] J. V. Burke, A. S. Lewis, and M. L. Overton. Two Numerical Methods for Optimizing Matrix Stability. *Linear Algebra Appl.*, 351/352:117–145, 2002.
- [BLO05] J. V. Burke, A. S. Lewis, and M. L. Overton. A Robust Gradient Sampling Algorithm for Nonsmooth, Nonconvex Optimization. *SIAM Journal on Optimization*, 15(3):751–779, 2005.
- [BMR14] E. Birgin, J. Martinez, and M. Raydan. Spectral Projected Gradient Methods: Review and Perspectives. *Journal of Statistical Software, Articles*, 60(3):1–21, 2014.
- [BO01] J. V. Burke and M. L. Overton. Variational Analysis of Non-Lipschitz Spectral Functions. *Math. Program.*, 90(2, Ser. A):317–351, 2001.
- [Cla75] F. H. Clarke. Generalized Gradients and Applications. *Trans. Amer. Math. Soc.*, 205:247–262, 1975.
- [Cla83] F. H. Clarke. *Optimization and Nonsmooth Analysis*. Wiley, New York, 1983. Reprinted by SIAM, Philadelphia, 1990.
- [CLR07] A. Crema, M. Loreto, and M. Raydan. Spectral Projected Subgradient with a Momentum Term for the Lagrangean Dual Approach. *Computers and Operations Research*, 34(10):3174–3186, 2007.
- [CMO17] F. E. Curtis, T. Mitchell, and M. L. Overton. A BFGS-SQP method for Nonsmooth, Nonconvex, Constrained Optimization and its Evaluation using Relative Minimization Profiles. *Optimization Methods and Software*, 32(1):148–181, 2017.
- [CO12] F. E. Curtis and M. L. Overton. A Sequential Quadratic Programming Algorithm for Nonconvex, Nonsmooth Constrained Optimization. *SIAM Journal on Optimization*, 22(2):474–500, 2012.
- [CQ13] F. E. Curtis and X. Que. An Adaptive Gradient Sampling Algorithm for Nonsmooth Optimization. *Optimization Methods and Software*, 28(6):1302–1324, 2013.
- [CQ15] F. E. Curtis and X. Que. A Quasi-Newton Algorithm for Nonconvex, Nonsmooth Optimization with Global Convergence Guarantees. *Mathematical Programming Computation*, 7(4):399–428, 2015.



- [CRZ17] F. E. Curtis, D. P. Robinson, and B. Zhou. Self-Correcting Variable-Metric Algorithms for Nonsmooth Optimization. Technical Report 17T-012, COR@L Laboratory, Department of ISE, Lehigh University, 2017.
- [DD18] D. Davis and D. Drusvyatskiy. Stochastic Model-Based Minimization of Weakly Convex Functions, March 2018. arXiv:1803.06523.
- [DDKL18] D. Davis, D. Drusvyatskiy, S. Kakade, and J. D. Lee. Stochastic Subgradient Method Converges on Tame Functions, April 2018. arXiv:1804.07795.
- [EM18] A. Estrada and I. M. Mitchell. Control Synthesis and Classification for Unicycle Dynamics using the Gradient and Value Sampling Particle Filters. In *Proceedings of the IFAC Conference on Analysis and Design of Hybrid Systems*, pages 108–114, July 2018. To appear.
- [Fle87] R. Fletcher. *Practical Methods of Optimization*. Wiley-Interscience, New York, NY, USA, Second edition, 1987.
- [Fle05] R. Fletcher. On the Barzilai-Borwein Method. In L. Qi, K. Teo, and X. Yang, editors, *Optimization and Control with Applications*, pages 235–256. Springer US, Boston, MA, USA, 2005.
- [GL17] J. Guo and A. S. Lewis. Nonsmooth Variants of Powell’s BFGS Convergence Theorem. *SIAM J. Optim.*, 2017. To appear.
- [Gol77] A. A. Goldstein. Optimization of Lipschitz Continuous Functions. *Math. Programming*, 13(1):14–22, 1977.
- [HN13] W. Hare and J. Nutini. A Derivative-Free Approximate Gradient Sampling Algorithm for Finite Minimax Problems. *Computational Optimization and Applications*, 56(1):1–38, Sep 2013.
- [HSS16] E. S. Helou, S. A. Santos, and L. E. A. Simões. On the Differentiability Check in Gradient Sampling Methods. *Optimization Methods and Software*, 31(5):983–1007, 2016.
- [HSS17] E. S. Helou, S. A. Santos, and L. E. A. Simões. On the Local Convergence Analysis of the Gradient Sampling Method for Finite Max-Functions. *Journal of Optimization Theory and Applications*, 175(1):137–157, 2017.
- [HU17] S. Hosseini and A. Uschmajew. A Riemannian Gradient Sampling Algorithm for Nonsmooth Optimization on Manifolds. *SIAM Journal on Optimization*, 27(1):173–189, 2017.
- [Kiw86] K. C. Kiwiel. A Method for Solving Certain Quadratic Programming Problems Arising in Nonsmooth Optimization. *IMA Journal of Numerical Analysis*, 6(2):137–152, 1986.
- [Kiw07] K. C. Kiwiel. Convergence of the Gradient Sampling Algorithm for Nonsmooth Nonconvex Optimization. *SIAM Journal on Optimization*, 18(2):379–388, 2007.
- [Kiw10] K. C. Kiwiel. A Nnderivative Version of the Gradient Sampling Algorithm for Nonsmooth Nonconvex Optimization. *SIAM Journal on Optimization*, 20(4):1983–1994, 2010.
- [LACR17] M. Loreto, H. Aponte, D. Cores, and M. Raydan. Nonsmooth Spectral Gradient Methods for Unconstrained Optimization. *EURO Journal on Computational Optimization*, 5(4):529–553, 2017.
- [Lew02] A. S. Lewis. Active Sets, Nonsmoothness, and Sensitivity. *SIAM Journal on Optimization*, 13(3):702–725, 2002.
- [Lin09] Q. Lin. *Sparsity and Nonconvex Nonsmooth Optimization*. PhD thesis, Department of Mathematics, University of Washington, 2009.

- [LMW16] J. Larson, M. Menickelly, and S. M. Wild. Manifold Sampling for  $\ell_1$  Nonconvex Optimization. *SIAM Journal on Optimization*, 26(4):2540–2563, 2016.
- [LO13] A. S. Lewis and M. L. Overton. Nonsmooth Optimization via Quasi-Newton Methods. *Math. Program.*, 141(1-2, Ser. A):135–163, 2013.
- [LOS00] 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.
- [MS05] R. Mifflin and C. Sagastizábal. A VU-Algorithm for Convex Minimization. *Mathematical Programming*, 104(2-3):583–608, 2005.
- [NS17] Y. Nesterov and V. Spokoiny. Random Gradient-Free Minimization of Convex Functions. *Found. Comput. Math.*, 17(2):527–566, 2017.
- [Ray93] M. Raydan. On the Barzilai and Borwein Choice of Steplength for the Gradient Method. *IMA Journal of Numerical Analysis*, 13(3):321–326, 1993.
- [Ray97] M. Raydan. The Barzilai and Borwein Gradient Method for the Large Scale Unconstrained Minimization Problem. *SIAM Journal on Optimization*, 7(1):26–33, 1997.
- [Roc82] R. T. Rockafellar. Lagrange Multipliers and Subderivatives of Optimal Value Functions in Nonlinear Programming. In D. C. Sorensen and R. J.-B. Wets, editors, *Mathematical Programming Study*, Mathematical Programming Studies, chapter 3, pages 28–66. North-Holland Publishing Company, Amsterdam, 1982.
- [Roc98] R. J.-B. Rockafellar, R. T. and Wets. *Variational analysis*, volume 317 of *Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences]*. Springer-Verlag, Berlin, 1998.
- [TLJL14] C.-M. Tang, S. Liu, J.-B. Jian, and J.-L. Li. A Feasible SQP-GS Algorithm for Nonconvex, Nonsmooth Constrained Optimization. *Numerical Algorithms*, 65(1):1–22, Jan 2014.
- [TM16] N. Traft and I. M. Mitchell. Improved Action and Path Synthesis using Gradient Sampling. In *Proceedings of the IEEE Conference on Decision and Control*, pages 6016–6023, December 2016.