



Globally convergent derivative-free methods in nonconvex optimization with and without noise

Pham Duy Khanh¹ · Boris S. Mordukhovich² · Dat Ba Tran²

Received: 25 June 2024 / Accepted: 29 June 2025

© Springer-Verlag GmbH Germany, part of Springer Nature and Mathematical Optimization Society 2025

Abstract

This paper addresses the study of nonconvex derivative-free optimization problems, where only information of either smooth objective functions or their noisy approximations is available. General derivative-free methods are proposed for minimizing differentiable (not necessarily convex) functions with globally Lipschitz continuous gradients, where the accuracy of approximate gradients is interacting with stepsizes and exact gradient values. Analysis in the noiseless case guarantees convergence of the gradient sequence to the origin as well as global convergence with constructive convergence rates of the sequence of iterates under the Kurdyka-Łojasiewicz property. In the noisy case, without any noise level information, the designed algorithms reach near-stationary points with providing estimates on the required number of iterations and function evaluations. Addressing functions with locally Lipschitzian gradients, two algorithms are introduced to handle the noiseless and noisy cases, respectively. The noiseless version is based on the standard backtracking linesearch and achieves fundamental convergence properties similarly to the global Lipschitzian case. The noisy version is based on a novel dynamic step linesearch and is shown to reach near-stationary points after a finite number of iterations when the Polyak-

B.S.Mordukhovich: Research of this author was partly supported by the US National Science Foundation under grants DMS-1808978 and DMS-2204519, by the Australian Research Council under grant DP-190100555, and by Project 111 of China under grant D21024. D. B. Tran: Research of this author was partly supported by the US National Science Foundation under grants DMS-1808978 and DMS-2204519. P. D. Khanh: This research was funded by the National Key Program for the Development of Mathematics in the period 2021–2030 under the National Science and Technology Project titled “Developing fundamental algorithms for finding optimal paths in 2.5 (terrain)- and 3-dimensional spaces” (Project Code: DTDLCN.05/25).

✉ Pham Duy Khanh
khanhpd@hcmue.edu.vn

Boris S. Mordukhovich
aa1086@wayne.edu

Dat Ba Tran
tranbadat@wayne.edu

¹ Group of Analysis and Applied Mathematics, Department of Mathematics, Ho Chi Minh City University of Education, Ho Chi Minh City, Vietnam

² Department of Mathematics, Wayne State University, Detroit, Michigan, USA

Łojasiewicz inequality is imposed. Numerical experiments are conducted on a diverse set of test problems to demonstrate more robustness of the newly proposed algorithms in comparison with other finite-difference-based schemes and some highly efficient, production-ready codes from the SciPy library. The experiments also demonstrate that the newly proposed methods can be integrated with acceleration techniques from the literature of smooth optimization while significantly enhancing numerical performance and outperforming current state-of-the-art derivative-free algorithms.

Keywords Derivative-free optimization · Nonconvex smooth objective functions · Finite differences · Black-box optimization · Noisy optimization · Zeroth-order optimization · Globally convergent algorithms

Mathematics Subject Classification 90C25 · 90C26 · 90C30 · 90C56

1 Introduction

This paper is devoted to the development of novel *derivative-free methods* of solving unconstrained optimization problems given in the form

$$\text{minimize } f(x) \quad \text{subject to } x \in \mathbb{R}^n, \quad (1.1)$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a continuously differentiable (\mathcal{C}^1 -smooth) function, not necessarily convex. In the context of derivative-free optimization, we assume that only information of either $f(x)$ (noiseless case) or its *noisy approximation* $\phi(x) = f(x) + \xi(x)$ (noisy case) is available, where $\xi : \mathbb{R}^n \rightarrow \mathbb{R}$ is the noise function bounded by a positive constant ξ_f . These problems have received much attention with a variety of methods being developed over the years [5, 17]. The major developments in this vein are provided by the *Nelder-Mead simplex method* [45, 51], *direct search methods* [31, 34], *conjugate direction method* [57], *trust-region methods* [18, 57], and *finite-difference-based methods* [6, 8, 52, 63, 65]. Applications of derivative-free optimization methods [2, 5, 17] have also gained a lot of interest since many efficient methods, including Nelder-Mead, Powell (a short name of Powell's conjugate direction method) [57], COBYLA [58], and L-BFGS-B, are implemented as production-ready codes in SciPy [66], a well-known Python library. More recently, numerous empirical results conducted by Shi et al. in [65] show that derivative-free optimization methods based on *finite differences* are accurate, efficient, and in some cases superior to other state-of-the-art derivative-free optimization methods developed in the literature. Meanwhile, extensive numerical comparisons in Berahas et al. [7], together with further analysis by Scheinberg [61], also tell us that the accuracy of gradients obtained from standard finite differences is significantly higher than from *randomized schemes* [21, 28, 29, 52]. These empirical results suggest that the methods using standard finite differences have much to be recommended, and that the research in this direction should be strongly encouraged.

When no error is present within the function evaluations, implementing finite difference approximations for gradient descent methods is rather simple because it is

possible to use a fixed, sufficiently small finite difference interval (referred to as *GD (fixed)* for the sake of brevity). However, dealing with noisy problems is more challenging since finding the optimal finite difference interval requires not only the *noise level* information but also the higher-order derivatives of the function that are often unavailable. Hence this topic attracts many studies, which develop finite-difference-based methods under different types of noise. Kelley et al. [16, 26, 38] proposed the *implicit filtering* algorithm based on a finite difference approach to deal with *noisy smooth box-constrained* optimization problems with the noise being decayed near local minimizers. Berahas et al. [6, 8] developed finite-difference-based linesearch methods for the minimization of *smooth functions* with *bounded noise*. The schemes to adapt the *finite difference intervals* were also studied by Gill et al. [25], Moré and Wild [49], and recently by Shi et al. [63, 64].

Motivations. Although methods of this type are often used in practice to solve derivative-free smooth problems with and without noise, there are still some significant concerns related to their theoretical and practical developments that should be addressed.

- *Analysis in the noiseless case:* In the noiseless case, due to the usage of a fixed finite difference interval, GD (fixed) methods do not obtain sufficient convergence properties compared to standard gradient descent methods. These properties include the *stationarity of accumulation points* and the *convergence* of the *sequence of iterates* to *nonisolated stationary points* under the *Kurdyka-Łojasiewicz* (KL) condition [3, 40–42, 44, 47], which is a rather mild regularity condition satisfied for the vast majority of objective functions in practice. Implicit filtering [38], which allows the finite difference interval to approach zero, is regarded as a more practical approach in the literature on derivative-free optimization with the stationarity of accumulation points guaranteed. However, the method has limitations in both practical and theoretical aspects as simply allowing the finite difference interval to approach zero without careful adaptive modifications, may result in an approximate gradient that is not even a descent direction. The versions presented in [39, Theorem 1] and [38, Theorem 7.1.1] require that the number of iterations with *failed linesearch* be finite, a condition not universally guaranteed. Similarly, the version in [53, Theorem 9.2] mandates that approximate gradients remain smaller than the finite difference interval, another assumption that may not always hold. From a broader perspective, the convergence of gradient descent with decreasing finite difference intervals can be derived from results on inexact gradient descent methods with decreasing error. However, its general applicability remains unclear. To the best of our knowledge, the most general related results with fundamental assumptions were presented in [10, Proposition 1] and [43, Theorem 3.3]. Both results necessitate that the error diminishes at a specific rate consistent with the stepsize. Furthermore, convergence rates are not analyzed in those works, which is a crucial aspect for achieving better numerical performance in methods of this type.
- *Dealing with small noise without any noise level information:* The practical implementations of finite-difference-based algorithms also face issues in this case since choosing sufficiently small finite difference intervals makes GD (fixed) methods

perform *poorly* due to the roundoff error, while using an adaptive scheme as in [25, 49, 63] becomes *inefficient* since the noise level is unknown. Although some methods for approximating the noise level may exist and be helpful in practice, their usage could significantly increase computational costs. This is particularly evident when the noise is not independent and identically distributed, which requires approximations for local noise levels to be conducted at every iteration.

- *Assumption on the gradient global Lipschitz continuity, i.e., the $\mathcal{C}_L^{1,1}$ property of objective functions:* This assumption seems to be omnipresent in derivative-free linesearch methods; see, e.g., [16, Theorem 2.1], [6, Assumption A1], and [8, Assumption 1.1]. For general derivative-free trust-region methods, Conn et al. [18] proved global convergence results in the case of smooth minimization problems while assuming the Lipschitz continuity of either the gradient or the Hessian; see [18, Assumptions 3.1 and 3.2]. Such properties were employed in the proximal point method adapted to derivative-free smooth optimization problems by Hare and Lucet [Assumption 1]. In [62], the class of smooth functions with Hölderian gradients, being larger than the class of $\mathcal{C}_L^{1,1}$ functions, was studied. However, we are not familiar with any efficient finite-difference-based method considering specifically the class of smooth functions with *locally Lipschitzian gradients*, i.e., the class of $\mathcal{C}^{1,1}$ functions, which is much broader than the class of $\mathcal{C}_L^{1,1}$ ones. This is in contrast to the exact versions of gradient descent methods that obtain various convergence properties including stationarity of accumulation points for the version with backtracking stepsizes addressing \mathcal{C}^1 -smooth functions [9, Proposition 1.2.1] and the global convergence for the version with sufficiently small stepsizes in the class of definable $\mathcal{C}^{1,1}$ functions [36]. This raises the need for the design and analysis of finite-difference-based methods concerning the class of $\mathcal{C}^{1,1}$ functions, which is the best we can hope in this context since the error bounds for finite differences are not available outside of this class.

Contributions. Having in mind the above discussions, we first address $\mathcal{C}_L^{1,1}$ optimization problems and propose the *derivative-free method with constant stepsize* (DFC). The method offers *generalizations* and *improvements* in both theoretical and practical aspects compared to GD (fixed) algorithms. The *generalizations* imply that other gradient approximation methods can be employed in DFC besides finite differences, provided that the approximation methods adhere to general conditions outlined in Definition 3.1. The *improvements* of DFC in comparison with GD (fixed) algorithms are as follows.

- In the *theoretical aspects*, DFC comes with a detailed convergence analysis, which is presented in Section 4 and contains the following.
 - In the *noiseless* case, under standard assumptions, our analysis establishes the *convergence to the origin* for the gradient sequence, *global convergence* of iterates under the KL property, and *constructive convergence rates* depending on the KL exponents. Note that none of these properties can be achieved by GD (fixed) algorithms, and while the implicit filtering algorithm [38] fulfills some them, certain nonstandard additional assumptions are further required.

- In the *noisy case*, the *finite convergence* of the sequence of iterates to a *near-stationary point* is established, along with estimates on the number of iterations and function evaluations needed to reach the near-stationary point. The construction of DFC and all of its convergence properties in the noisy case do *not require* the knowledge of *noise levels*, although it does require some mild conditions for initialization.
- In the *practical aspects*, DFC achieves at least similar, or even better, numerical performance in comparison with GD (fixed) methods in the *noiseless* case being more efficient in the presence of small noise with an *unknown noise level*. These numerical results are presented in Section 6.

Note that the main feature in the algorithmic constructions of DFC, which allows us to achieve the above goals, is the *adaptivity* of the finite difference interval. Contrarily to using a fixed small interval as in GD (fixed) methods, we start with a *much larger* finite difference interval and decrease it along the sequence of iterates if a descent condition is not satisfied. The finite difference interval in DFC also interacts with the approximate Lipschitz constant (or equivalently, the stepsize), which creates more robustness for the algorithm. This interaction distinguishes DFC from the methods in [25, 49, 63, 64], where the finite difference interval is constructed independently from the stepsize. By adopting this approach, we are able to theoretically derive the *fundamental convergence properties* of DFC for both noiseless and noisy functions. Practically this approach helps DFC *avoiding roundoff errors* as much as possible to ensure the quality of the gradient approximation, which leads us to the numerical performance highlighted above.

Next we address the class of $\mathcal{C}^{1,1}$ functions. Due to the complex structure of functions in this class, we introduce two different algorithms as follows.

- DFB: *Derivative-free method with backtracking linesearch* to deal with noiseless problems and problems with small noise. This algorithm is inspired by DFC, with the primary difference that a backtracking linesearch step is performed in *each iteration*, similarly to the standard approach of gradient descent methods when dealing with $\mathcal{C}^{1,1}$ functions. The analysis is conducted in the noiseless case and establishes the *stationarity of the accumulation point*, *global convergence* under the KL property, and *constructive convergence rates* depending on the KL exponents.
- DFD: *Derivative-free method with dynamic step linesearch* to deal with problems with large noise and known noise level. To be more specific, in each iteration DFD uses a dynamic step linesearch to approximate the local Lipschitz constant of the gradient in a region around the current iterate. The approximate Lipschitz constant is then used for determining both the stepsize and the finite difference interval. By employing this approach, DFD exhibits more favorable numerical behavior compared to other finite difference schemes [26, 52, 63] as demonstrated in Figure 2. This example also shows that the standard backtracking linesearch does not work well for functions of $\mathcal{C}^{1,1}$ class with *large noise*, which is the main motivation for us to implement the dynamic step linesearch in DFD. The global analysis of DFD demonstrates that the algorithm always makes progress whenever the gradient at the current iterate is not near the origin. It is also established that the sequence of iterates finds a near-stationary point after a *finite number of iterations*,

with a constructive estimate given when the *Polyak-Łojasiewicz inequality* (cf. [55] and [47]) is satisfied.

To demonstrate the practical aspects of our study, *extensive numerical experiments* on synthetic problems with and without noise are conducted in Section 6. The results when the noise is small show that DFC and DFB methods do improve the performance of GD (fixed) methods as well as the performance of the implicit filtering (IMFIL) algorithm [26] and random gradient-free (RG) algorithm [52]. When the noise is large, our DFD demonstrates its numerical reliability in comparison with SciPy [66] production-ready codes, including Powell, COBYLA and L-BFGS-B algorithms.

It is also demonstrated that, similar to standard gradient descent methods, the numerical performance of DFC can be significantly improved by incorporating additional steps such as those based on either quasi-Newton techniques [53], or Polyak momentum [55]. The experiments show that these variants not only enhance the basic version of DFC but also outperform the state-of-the-art Powell method from the Scipy library in most cases. Given the extensive development of gradient-based methods in the optimization literature, this phenomenon highlights not only the potential of the algorithms but also opens new research directions for developing DFC variants that are efficient in diverse scenarios and come with rigorous theoretical guarantees.

Related Works. The adaptivity of the finite difference interval to ensure the quality of the approximate gradient is a main feature employed in many finite-difference-based methods. Cartis and Scheinberg [14] analyzed a general linesearch algorithm for smooth functions without noise under the major condition that the gradient estimates are sufficiently accurate with a certain probability. This analysis is then extended in Berahas et al. [6] to the case where the function values are noisy. The practical schemes to choose the finite difference intervals adaptively based on testing ratios were also studied by Gill et al. [25], Shi et al. [63, 64], and heuristically by Moré and Wild [49]. The adaptivity in the selection of the finite difference interval is also related to the dynamic accuracy of gradient approximations, which is considered for adaptive regularization algorithms without noise in [13, 30] and with noise in [12, 15]. Recently, [20] also employed finite differences for implementing regularized Newton methods, which adaptively adjust the finite difference interval and link it with the cubic regularization parameter.

Among the aforementioned publications, [14] and [6] are the most related to our DFC development for $\mathcal{C}_L^{1,1}$ functions. These results are discussed in more detail in Remark 4.4 and Remark 4.11.

Another type of optimization methods commonly used in the derivative-free setting is derived from model-based trust-region algorithms as presented in [17, 18] and the references therein. There are significant differences arising from the fundamental ideas behind the construction of trust-region methods and the linesearch algorithms presented in this work. The former methods construct local models to approximate the objective function at each iteration and restrict the updates to a region where this approximation is reliable. In [17, Chapter 10], the size of the trust region is adaptively adjusted based on how well the model predicts the actual improvement in the objective function while attempting to balance both exploration and exploitation of iterates. This also highlights that both the accuracy of the model and the stepsize of the iterate

are defined by the trust-region. In contrast, our approach uses finite differences to approximate the gradients and employs them as the direction for updating the iterate without any restriction on the stepsize. Additionally, the accuracy of the gradient and stepsize are determined separately. Both the direction and stepsize of our method can be flexibly adjusted by incorporating acceleration techniques, further emphasizing its advantages.

Organization. The rest of the paper is organized as follows. Section 2 presents some basic definitions and preliminaries used throughout the entire paper. Section 3 examines two types of gradient approximations that include finite differences. The main parts of our work, concerning the design and convergence properties of general derivative-free methods under the global and local Lipschitz continuity of the gradient, are given in Section 4 and Section 5, respectively. Numerical experiments, which compare the efficiency of the proposed methods with other derivative-free methods for both noisy and noiseless functions, are conducted in Section 6. Concluding remarks on the main contributions of this paper together with some open questions and perspectives of our future research are presented in Section 7.

2 Preliminaries

First we recall some basic notions and notation frequently used in the paper. All our considerations are given in the space \mathbb{R}^n with the Euclidean norm $\|\cdot\|$. For any $i = 1, \dots, n$, let e_i denote the i^{th} basic vector in \mathbb{R}^n . As always, $\mathbb{N} := \{1, 2, \dots\}$ signifies the collection of natural numbers. For any $x \in \mathbb{R}^n$ and $\varepsilon > 0$, let $\mathbb{B}(x, \varepsilon)$ and $\overline{\mathbb{B}}(x, \varepsilon)$ stand for the open and closed balls centered at x with radius ε , respectively. When $x = 0$, these balls are denoted simply by $\varepsilon\mathbb{B}$ and $\varepsilon\overline{\mathbb{B}}$.

Recall that a mapping $G: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is *Lipschitz continuous on a subset D of \mathbb{R}^n* if there exists a constant $L > 0$ such that we have

$$\|G(x) - G(y)\| \leq L \|x - y\| \text{ for all } x, y \in D.$$

If $D = \mathbb{R}^n$, the mapping G is said to be *globally Lipschitz continuous*. The *local Lipschitz continuity* of G on \mathbb{R}^n is understood as the Lipschitz continuity of this mapping on every compact subset of \mathbb{R}^n . The latter is equivalent to saying that for any $x \in \mathbb{R}^n$ there is a neighborhood U of x such that G is Lipschitz continuous on U . In what follows, we denote by $\mathcal{C}^{1,1}$ the class of \mathcal{C}^1 -smooth mappings that have a *locally Lipschitz continuous gradient* on \mathbb{R}^n and by $\mathcal{C}_L^{1,1}$ the class of \mathcal{C}^1 -smooth mappings that have a *globally Lipschitz continuous gradient with the constant $L > 0$* (i.e., L -Lipschitz continuous) on the entire space.

Our convergence analysis of the numerical algorithms developed in the subsequent sections largely exploits the following important results and notions. The first result taken from [35, Lemma A.11] presents a simple albeit very useful property of real-valued functions with Lipschitz continuous gradients.

Lemma 2.1 *Let $f: \mathbb{R}^n \rightarrow \mathbb{R}$, let $x, y \in \mathbb{R}^n$, and let $L > 0$. If f is differentiable on the line segment $[x, y]$ with its derivative being L -Lipschitz continuous on this*

segment, then

$$|f(y) - f(x) - \langle \nabla f(x), y - x \rangle| \leq \frac{L}{2} \|y - x\|^2. \quad (2.1)$$

The second lemma established in [40, Section 3] is crucial in the convergence analysis of the general linesearch methods developed in this paper.

Lemma 2.2 *Let $\{x^k\}$ and $\{d^k\}$ be sequences in \mathbb{R}^n satisfying the condition*

$$\sum_{k=1}^{\infty} \|x^{k+1} - x^k\| \cdot \|d^k\| < \infty. \quad (2.2)$$

If \bar{x} is an accumulation point of $\{x^k\}$ and if the origin is an accumulation point of $\{d^k\}$, then there exists an infinite set $J \subset \mathbb{N}$ such that

$$x^k \xrightarrow{J} \bar{x} \text{ and } d^k \xrightarrow{J} 0. \quad (2.3)$$

Next we recall the classical results from [22, Section 8.3.1] that describe important properties of accumulation points generated by a sequence satisfying the limit condition introduced by Ostrowski [54].

Lemma 2.3 *Let $\{x^k\} \subset \mathbb{R}^n$ be a sequence satisfying the Ostrowski condition*

$$\lim_{k \rightarrow \infty} \|x^{k+1} - x^k\| = 0. \quad (2.4)$$

Then the following assertions are fulfilled:

- (i) *If $\{x^k\}$ is bounded, then the set of accumulation points of $\{x^k\}$ is nonempty, compact, and connected in \mathbb{R}^n .*
- (ii) *If $\{x^k\}$ has an isolated accumulation point \bar{x} , then this sequence converges to \bar{x} .*

The version of the fundamental Kurdyka-Łojasiewicz (KL) property formulated below is taken from Absil et al. [1, Theorem 3.4].

Definition 2.4 Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a differentiable function. We say that f satisfies the KL property at $\bar{x} \in \mathbb{R}^n$ if there exist a number $\eta > 0$, a neighborhood U of \bar{x} , and a nondecreasing function $\psi : (0, \eta) \rightarrow (0, \infty)$ such that the function $1/\psi$ is integrable over $(0, \eta)$ and we have

$$\|\nabla f(x)\| \geq \psi(f(x) - f(\bar{x})) \text{ for all } x \in U \text{ with } f(\bar{x}) < f(x) < f(\bar{x}) + \eta. \quad (2.5)$$

Remark 2.5 If f satisfies the KL property at \bar{x} with a neighborhood U , it is clear that the same property holds for any $x \in U$ where $f(x) = f(\bar{x})$. It has been realized that the KL property is satisfied in broad settings. In particular, it holds at every *nonstationary point*

of f ; see [3, Lemma 2.1 and Remark 3.2(b)]. Furthermore, it is proved in the seminal paper by Łojasiewicz [47] that any analytic function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ satisfies the KL property at every point \bar{x} with $\psi(t) = Mt^q$ for some $q \in [0, 1)$. As demonstrated in [40, Section 2], the KL property formulated in Attouch et al. [3] is stronger than the one in Definition 2.4. Typical smooth functions that satisfy the KL property from [3], and hence the one from Definition 2.4, are smooth *semialgebraic* functions and also those from the more general class of functions known as *definable in o-minimal structures*; see [3, 4, 44]. The latter property is fulfilled, e.g., in important models arising in deep neural networks, low-rank matrix recovery, principal component analysis, and matrix completion as discussed in [11, Section 6.2].

Next we present, based on [1], some descent-type conditions ensuring the global convergence of iterates for smooth functions that satisfy the KL property.

Proposition 2.6 *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a C^1 -smooth function, and let the sequence of iterations $\{x^k\} \subset \mathbb{R}^n$ satisfy the following conditions:*

(H1) (primary descent condition). *There exists $\sigma > 0$ such that for sufficiently large $k \in \mathbb{N}$, we have*

$$f(x^k) - f(x^{k+1}) \geq \sigma \|\nabla f(x^k)\| \cdot \|x^{k+1} - x^k\|.$$

(H2) (complementary descent condition). *For sufficiently large $k \in \mathbb{N}$, we have*

$$[f(x^{k+1}) = f(x^k)] \implies [x^{k+1} = x^k].$$

If \bar{x} is an accumulation point of $\{x^k\}$ and f satisfies the KL property at \bar{x} , then $x^k \rightarrow \bar{x}$ as $k \rightarrow \infty$.

When the sequence under consideration is generated by a linesearch method and satisfies some conditions stronger than (H1) and (H2) in Proposition 2.6, its convergence rates are established in [40, Proposition 2.4] under the KL property with $\psi(t) = Mt^q$ as given below.

Proposition 2.7 *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a C^1 -smooth function, and let the sequences $\{x^k\} \subset \mathbb{R}^n$, $\{\tau_k\} \subset [0, \infty)$, $\{d^k\} \subset \mathbb{R}^n$ satisfy the iterative condition $x^{k+1} = x^k + \tau_k d^k$ for all $k \in \mathbb{N}$. Assume that for sufficiently large $k \in \mathbb{N}$, we have $x^{k+1} \neq x^k$ together with the estimates*

$$f(x^k) - f(x^{k+1}) \geq \beta \tau_k \|d^k\|^2 \quad \text{and} \quad \|\nabla f(x^k)\| \leq \alpha \|d^k\|, \quad (2.6)$$

where α, β are some positive constants. Suppose in addition that the sequence $\{\tau_k\}$ is bounded away from 0 (i.e., there exists some $\bar{\tau} > 0$ such that $\tau_k \geq \bar{\tau}$ for sufficiently large $k \in \mathbb{N}$), that \bar{x} is an accumulation point of $\{x^k\}$, and that f satisfies the KL property at \bar{x} with $\psi(t) = Mt^q$ for some $M > 0$ and $q \in [1/2, 1)$. Then the following convergence rates are guaranteed:

- (i) If $q = 1/2$, then the sequence $\{x^k\}$ converges linearly to \bar{x} .
(ii) If $q \in (1/2, 1)$, then we have the estimate

$$\|x^k - \bar{x}\| = \mathcal{O}\left(k^{-\frac{1-q}{2q-1}}\right).$$

Remark 2.8 Observe that the two conditions in (2.6) together with the boundedness away from 0 of $\{\tau_k\}$ yield assumptions (H1), (H2) in Proposition 2.6. Indeed, (H1) is verified by the following inequalities:

$$\begin{aligned} f(x^k) - f(x^{k+1}) &\geq \beta \tau_k \|d^k\|^2 = \beta \|\tau_k d^k\| \cdot \|d^k\| \\ &\geq \frac{\beta}{\alpha} \|x^{k+1} - x^k\| \cdot \|\nabla f(x^k)\|. \end{aligned}$$

In addition, since $\{\tau_k\}$ is bounded away from 0, there exists $\bar{\tau} > 0$ such that $\tau_k \geq \bar{\tau}$ for sufficiently large $k \in \mathbb{N}$. Then for such k , the condition $f(x^{k+1}) = f(x^k)$ implies that $d^k = 0$ by the first inequality in (2.6), and hence $x^{k+1} = x^k$ by the iterative procedure $x^{k+1} = x^k + \tau_k d^k$, which therefore verifies (H2).

3 Global and local approximations of gradients

This section is devoted to analyzing several methods for approximating gradients of a smooth function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ by using only information about the function values that frequently appears in derivative-free optimization. Methods of this type include, in particular, finite differences [53, Section 9], the Gupal estimation [33], and gradient estimation via linear interpolation [7]. We construct two types of approximations that cover all these methods.

Definition 3.1 Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a C^1 -smooth function. A mapping $\mathcal{G} : \mathbb{R}^n \times (0, \infty) \rightarrow \mathbb{R}^n$ is:

- (i) A *global approximation* of ∇f if there is a constant $C > 0$ such that

$$\|\mathcal{G}(x, \delta) - \nabla f(x)\| \leq C\delta \quad \text{for any } (x, \delta) \in \mathbb{R}^n \times (0, \infty). \quad (3.1)$$

- (ii) A *local approximation* of ∇f if for any bounded set $\Omega \subset \mathbb{R}^n$ and any $\Delta > 0$, there is $C > 0$ with

$$\|\mathcal{G}(x, \delta) - \nabla f(x)\| \leq C\delta \quad \text{for any } (x, \delta) \in \Omega \times (0, \Delta]. \quad (3.2)$$

Remark 3.2 We have the following observations related to Definition 3.1:

- (i) If \mathcal{G} is a global approximation of ∇f , then it is also a local approximation of ∇f .

(ii) Assume that \mathcal{G} is a local approximation of ∇f and that $x \in \mathbb{R}^n$. Then we deduce from (3.2) with $\Omega = \{x\}$ and any $\Delta > 0$ the condition

$$\limsup_{\delta \downarrow 0} \frac{\|\mathcal{G}(x, \delta) - \nabla f(x)\|}{\delta} < \infty. \quad (3.3)$$

Next we formulate the two standard types of finite differences taken from [53, Section 9], which serve as typical examples of the approximations in Definition 3.1.

• *Forward finite difference:*

$$\mathcal{G}(x, \delta) := \frac{1}{\delta} \sum_{i=1}^n (f(x + \delta e_i) - f(x)) e_i \text{ for any } (x, \delta) \in \mathbb{R}^n \times (0, \infty). \quad (3.4)$$

• *Central finite difference:*

$$\mathcal{G}(x, \delta) := \frac{1}{2\delta} \sum_{i=1}^n (f(x + \delta e_i) - f(x - \delta e_i)) e_i \text{ for any } (x, \delta) \in \mathbb{R}^n \times (0, \infty). \quad (3.5)$$

In the constructions above, the positive number δ is called the *finite difference interval*.

Remark 3.3 Let us now recall some results on the error bounds for the two types of finite differences that are mentioned above.

(i) The global error bound for the forward finite difference (see, e.g., [7, Theorem 2.1] and [53, Section 8]) shows that it is a global approximation of ∇f when $f \in \mathcal{C}_L^{1,1}$. The local error bound for the forward finite difference is also given in [53, Exercise 9.13].

(ii) On the other hand, the global error bound for the central finite difference (see, e.g., [7, Theorem 2.2] and [53, Lemma 9.1]) requires that f is twice continuously differentiable with a Lipschitz continuous Hessian, which is a rather restrictive assumption.

For completeness, we present a short proof showing that both types of finite differences are global approximations of ∇f when $f \in \mathcal{C}_L^{1,1}$ and are local approximations of ∇f when $f \in \mathcal{C}^{1,1}$.

Proposition 3.4 *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a \mathcal{C}^1 -smooth function. Then the following hold:*

(i) *Given $x \in \mathbb{R}^n$ and $\delta > 0$, if the gradient ∇f is Lipschitz continuous on $\mathbb{B}(x, \delta)$ with the constant $L > 0$, then both forward finite difference (3.4) and central finite difference (3.5) satisfy the estimate*

$$\|\mathcal{G}(x, \delta) - \nabla f(x)\| \leq \frac{L\sqrt{n}\delta}{2}. \quad (3.6)$$

(ii) *If the gradient ∇f is globally Lipschitz continuous with the constant $L > 0$, then both forward finite difference (3.4) and central finite difference (3.5) satisfy the*

estimate

$$\|\mathcal{G}(x, \delta) - \nabla f(x)\| \leq \frac{L\sqrt{n}\delta}{2} \text{ for any } (x, \delta) \in \mathbb{R}^n \times (0, \infty). \quad (3.7)$$

(iii) If the gradient ∇f is locally Lipschitz continuous, then for any bounded set $\Omega \subset \mathbb{R}^n$ and for any $\Delta > 0$, there exists a positive number L such that both forward finite difference (3.4) and central finite difference (3.5) satisfy the estimate

$$\|\mathcal{G}(x, \delta) - \nabla f(x)\| \leq \frac{L\sqrt{n}\delta}{2} \text{ for any } (x, \delta) \in \Omega \times (0, \Delta]. \quad (3.8)$$

Proof We begin with verifying (i) for each type of the aforementioned finite differences and then employ (i) to justify (ii) and (iii) for both types.

(i) Take any $x \in \mathbb{R}^n$, $\delta > 0$ and assume that ∇f is Lipschitz continuous on $\overline{\mathbb{B}}(x, \delta)$ with the constant $L > 0$. Consider first the case where \mathcal{G} is given by the forward finite difference (3.4). Then for any $i = 1, \dots, n$, we get by employing Lemma 2.1 that

$$|f(x + \delta e_i) - f(x) - \langle \nabla f(x), x + \delta e_i - x \rangle| \leq \frac{L}{2} \|x + \delta e_i - x\|^2 = \frac{L\delta^2}{2}, \quad (3.9)$$

which is clearly equivalent to

$$\left| \frac{1}{\delta} (f(x + \delta e_i) - f(x)) - \frac{\partial f}{\partial x_i}(x) \right| \leq \frac{L\delta}{2}.$$

Since the latter inequality holds for all $i = 1, \dots, n$, we deduce that

$$\|\mathcal{G}(x, \delta) - \nabla f(x)\| = \sqrt{\sum_{i=1}^n \left(\frac{1}{\delta} (f(x + \delta e_i) - f(x)) - \frac{\partial f}{\partial x_i}(x) \right)^2} \leq \frac{L\sqrt{n}\delta}{2},$$

which therefore verifies estimate (3.6).

Assume now that \mathcal{G} is given by the central finite difference (3.5). Employing Lemma 2.1 gives us for any $i = 1, \dots, n$ the two estimates

$$\begin{aligned} |f(x + \delta e_i) - f(x) - \langle \nabla f(x), (x + \delta e_i) - x \rangle| &\leq \frac{L\delta^2}{2}, \\ |f(x) - f(x - \delta e_i) - \langle \nabla f(x), x - (x - \delta e_i) \rangle| &\leq \frac{L\delta^2}{2}. \end{aligned}$$

Summing up the above estimates and using the triangle inequality, we deduce that

$$|f(x + \delta e_i) - f(x - \delta e_i) - 2 \langle \nabla f(x), \delta e_i \rangle| \leq L\delta^2,$$

which implies in turn the conditions

$$\left| \frac{1}{2\delta} (f(x + \delta e_i) - f(x - \delta e_i)) - \frac{\partial f}{\partial x_i}(x) \right| \leq \frac{L\delta}{2}$$

for all $i = 1, \dots, n$. Therefore, we get

$$\|\mathcal{G}(x, \delta) - \nabla f(x)\| = \sqrt{\sum_{i=1}^n \left(\frac{1}{2\delta} (f(x + \delta e_i) - f(x - \delta e_i)) - \frac{\partial f}{\partial x_i}(x) \right)^2} \leq \frac{L\sqrt{n}\delta}{2},$$

which brings us to (3.6) and thus justifies (i).

Assertion (ii) follows directly from (i). To verify (iii), pick some $\Delta > 0$ and a bounded set $\Omega \subset \mathbb{R}^n$, and then find $r > 0$ such that $\Omega \subset r\mathbb{B}$. Defining $\Theta := (r + \Delta)\mathbb{B}$, it is clear that Θ is compact. Since ∇f is locally Lipschitzian, it is Lipschitz continuous on Θ with some constant $L > 0$. Taking any $(x, \delta) \in \Omega \times (0, \Delta]$, we get that $\mathbb{B}(x, \delta) \subset \Theta$, and thus ∇f is Lipschitz continuous on $\mathbb{B}(x, \delta)$ with the same constant L . Employing finally (i) justifies assertion (iii). \square

The following example shows that when the local Lipschitz continuity of ∇f is replaced by merely the continuity of ∇f , the finite differences may not be a local approximation of ∇f .

Example 1 Define the univariate real-valued function f by

$$f(x) := \begin{cases} \frac{2}{3}\sqrt{x^3} & \text{if } x \geq 0, \\ -\frac{2}{3}\sqrt{-x^3} & \text{if } x < 0. \end{cases}$$

The derivative of f is calculated by

$$\nabla f(x) = \begin{cases} \sqrt{x} & \text{if } x \geq 0, \\ \sqrt{-x} & \text{if } x < 0 \end{cases}$$

being clearly continuous on \mathbb{R} while not Lipschitz continuous around 0. If we suppose that $\mathcal{G}(x, \delta)$ is the forward finite difference approximation of $\nabla f(x)$ from (3.4), we get that

$$\mathcal{G}(0, \delta) = \frac{f(\delta) - f(0)}{\delta} = \frac{\frac{2}{3}\sqrt{\delta^3}}{\delta} = \frac{2\sqrt{\delta}}{3} \text{ for all } \delta > 0,$$

which implies that $\mathcal{G}(0, \delta)/\delta \rightarrow \infty$ as $\delta \downarrow 0$. It follows from (3.3) that $\mathcal{G}(x, \delta)$ is not a local approximation of the derivative ∇f . Supposing now that $\mathcal{G}(x, \delta)$ is the central finite difference approximation of $\nabla f(x)$, we deduce from (3.5) the expression

$$\mathcal{G}(0, \delta) = \frac{f(\delta) - f(-\delta)}{2\delta} = \frac{4\sqrt{\delta}}{3} \text{ for all } \delta > 0,$$

which also tells us that $\mathcal{G}(x, \delta)$ is not a local approximation of ∇f .

4 General derivative-free methods for $\mathcal{C}_L^{1,1}$ functions

This section addresses the optimization problem (1.1) when $f \in \mathcal{C}_L^{1,1}$ for some $L > 0$. By employing gradient approximation methods that satisfy the global error bound (3.1), we propose the general *derivative-free method with constant stepsize* (DFC) to solve this problem for both noiseless and noisy cases, providing its convergence analysis. The DFC algorithm is described as follows.

4.1 Algorithm Construction

Algorithm 1 (DFC).

Step 0. Choose a global approximation \mathcal{G} of ∇f under condition (3.1). Select an initial point $x^1 \in \mathbb{R}^n$, an initial sampling radius $\delta_1 > 0$, a constant $C_1 > 0$, a reduction factor $\theta \in (0, 1)$, and scaling factors $\mu > 2$, $\eta > 1$, $\kappa > 0$. Set $k := 1$.

Step 1 (approximate gradient). Find g^k and the smallest nonnegative integer i_k such that

$$g^k = \mathcal{G}(x^k, \theta^{i_k} \delta_k) \text{ and } \|g^k\| > \mu C_k \theta^{i_k} \delta_k.$$

Then set $\delta_{k+1} := \theta^{i_k} \delta_k$.

Step 2 (update). If $f\left(x^k - \frac{\kappa}{C_k} g^k\right) \leq f(x^k) - \frac{\kappa(\mu - 2)}{2C_k \mu} \|g^k\|^2$, then $x^{k+1} := x^k - \frac{\kappa}{C_k} g^k$ and $C_{k+1} := C_k$. Otherwise, $x^{k+1} := x^k$ and $C_{k+1} := \eta C_k$.

Remark 4.1 Let us present some observations concerning Algorithm 1. The first observation clarifies the existence of g^k and i_k in Step 1. Observation (ii) explains the iteration updates in Step 2 while observation (iii) interprets the term “constant step-size” in the name of our method.

(i) The procedure of finding g^k and i_k that satisfy Step 1 can be given as follows. Set $i_k := 0$ and

$$g^k := \mathcal{G}(x^k, \theta^{i_k} \delta_k). \quad (4.1)$$

While $\|g^k\| \leq \mu C_k \theta^{i_k} \delta_k$, increase i_k by 1 and recalculate g^k under (4.1). When $\nabla f(x^k) \neq 0$, the existence of g^k and i_k in Step 1 is guaranteed. Indeed, otherwise we get a sequence $\{g_i^k\}$ with

$$g_i^k = \mathcal{G}(x^k, \theta^i \delta_k) \text{ and } \|g_i^k\| \leq \mu \theta^i \delta_k \text{ for all } i \in \mathbb{N}. \quad (4.2)$$

Since $\theta \in (0, 1)$, the latter means that $g_i^k \rightarrow 0$ as $i \rightarrow \infty$. Remembering that \mathcal{G} is a global approximation of ∇f , we get for $C > 0$ given in (3.1) that

$$\|g_i^k - \nabla f(x^k)\| \leq C\theta^i \delta_k \quad \text{whenever } i \in \mathbb{N}.$$

Letting $i \rightarrow \infty$ with taking into account that $g_i^k \rightarrow 0$, the latter inequality implies that $\nabla f(x^k) = 0$, which is a contradiction.

(ii) The condition $f(x^k - C_k^{-1} \kappa g^k) \leq f(x^k) - \frac{\kappa(\mu - 2)}{2C_k \mu} \|g^k\|^2$ determines whether C_k is a good approximation for C in the sense that the objective function f is sufficiently decreasing when the iterate moves from x^k to $x^{k+1} := x^k - C_k^{-1} \kappa g^k$. If this condition fails, we increase C_k by setting $C_{k+1} := \eta C_k$ to get a better approximation for C and stagnate the iterative sequence by setting $x^{k+1} := x^k$.

(iii) It will be shown in Proposition 4.2 that there exists a positive number \bar{C} such that $C_k = \bar{C}$ for sufficiently large $k \in \mathbb{N}$, which also implies that $x^{k+1} = x^k - \kappa \bar{C}^{-1} g^k$ for such k . This explains the term “constant stepsize” in the name of our algorithm.

(iv) The constant κ is a positive scaling factor that can be chosen arbitrarily. However, to ensure a good performance when finite difference approximations are used, κ is usually chosen as $\frac{\sqrt{n}}{2}$. In this case, by defining $L_k := \frac{C_k}{\kappa}$, the stepsize $\frac{\kappa}{C_k}$, as introduced in Step 2 above, becomes $\frac{1}{L_k}$, which is near the optimal stepsize in gradient descent if L_k is sufficiently close to the Lipschitz constant L of ∇f .

On the other hand, the constant $\mu > 2$ represents the relative error threshold for gradient approximation. This is evident from the fact that if $C_k \geq C$, it follows from Step 1 and estimate (3.1) that we have $\|g^k - \nabla f(x^k)\| \leq \frac{1}{\mu} \|g^k\|$.

4.2 Analysis for noiseless functions

In this subsection, we derive convergence properties of DFC in Algorithm 2 for noiseless functions, i.e., when $f(x)$ is available for all $x \in \mathbb{R}^n$. Our analysis begins with a crucial result showing that the tail of the sequence $\{C_k\}$ generated by Algorithm 1 is constant.

Proposition 4.2 *Let $\{C_k\}$ be the sequence generated by Algorithm 1. Assume that $\nabla f(x^k) \neq 0$ for all $k \in \mathbb{N}$. Then there exists a number $N \in \mathbb{N}$ such that $C_{k+1} = C_k$ whenever $k \geq N$.*

Proof Since \mathcal{G} is a global approximation of ∇f under condition (3.1), there exists $C > 0$ such that

$$\|\mathcal{G}(x, \delta) - \nabla f(x)\| \leq C\delta \quad \text{for all } (x, \delta) \in \mathbb{R}^n \times (0, \infty). \quad (4.3)$$

By the imposed assumption, we find $L > 0$ such that ∇f is Lipschitz continuous with the constant L on \mathbb{R}^n . Arguing by contradiction, suppose that the number N asserted in the proposition does not exist. By Step 2 of Algorithm 1, this implies that $C_{k+1} = \eta C_k$ for infinitely many $k \in \mathbb{N}$, and hence $C_k \rightarrow \infty$ as $k \rightarrow \infty$. Therefore,

there exists a number $K \in \mathbb{N}$ such that $C_{K+1} = \eta C_K$ and $C_K > \max \{C, L\kappa\}$. Using Step 2 of Algorithm 1 together with the update $C_{K+1} = \eta C_K$, we deduce that

$$f\left(x^K - \frac{\kappa}{C_K} g^K\right) > f(x^K) - \frac{\kappa(\mu - 2)}{2C_K\mu} \|g^K\|^2. \quad (4.4)$$

Combining $g^K = \mathcal{G}(x^K, \delta_{K+1})$ and $\|g^K\| \geq \mu C_K \delta_{K+1}$ from Step 1 of Algorithm 1 with (4.3) and $C_K > C$ as above, we get the relationships

$$\begin{aligned} \|g^K - \nabla f(x^K)\| &= \|\mathcal{G}(x^K, \delta_{K+1}) - \nabla f(x^K)\| \\ &\leq C\delta_{K+1} \leq C_K\delta_{K+1} \leq \mu^{-1} \|g^K\|. \end{aligned}$$

By the Cauchy-Schwarz inequality, the latter tells us that

$$\begin{aligned} \langle \nabla f(x^K), g^K \rangle &= \langle \nabla f(x^K) - g^K, g^K \rangle + \|g^K\|^2 \\ &\geq -\|\nabla f(x^K) - g^K\| \cdot \|g^K\| + \|g^K\|^2 \\ &\geq (1 - \mu^{-1}) \|g^K\|^2. \end{aligned}$$

Combining this with Lemma 2.1 and taking into account the global Lipschitz continuity of ∇f with the constant L as well as the condition $C_K > L\kappa$ as above, we get that

$$\begin{aligned} f\left(x^K - \frac{\kappa}{C_K} g^K\right) - f(x^K) &\leq -\frac{\kappa}{C_K} \langle \nabla f(x^K), g^K \rangle + \frac{L}{2} \left\| \frac{\kappa}{C_K} g^K \right\|^2 \\ &\leq -\frac{\kappa}{C_K} \left(1 - \frac{1}{\mu}\right) \|g^K\|^2 + \frac{\kappa}{2C_K} \|g^K\|^2 \\ &= -\frac{\kappa}{C_K} \|g^K\|^2 \left(\frac{1}{2} - \frac{1}{\mu}\right) = -\frac{\kappa(\mu - 2)}{2C_K\mu} \|g^K\|^2, \end{aligned}$$

which clearly contradicts (4.4) and thus completes the proof of the proposition. \square

Now we are ready to establish the convergence properties of Algorithm 1 in the noiseless case.

Theorem 4.3 *Let $\{x^k\}$ be the sequence generated by Algorithm 1 and assume that $\nabla f(x^k) \neq 0$ for all $k \in \mathbb{N}$. Then either $f(x^k) \rightarrow -\infty$ as $k \rightarrow \infty$, or we have the assertions:*

- (i) *The gradient sequence $\{\nabla f(x^k)\}$ converges to 0 as $k \rightarrow \infty$.*
- (ii) *If f satisfies the KL property at some accumulation point \bar{x} of $\{x^k\}$, then $x^k \rightarrow \bar{x}$ as $k \rightarrow \infty$.*
- (iii) *If f satisfies the KL property at some accumulation point \bar{x} of $\{x^k\}$ with $\psi(t) = Mt^q$ for $M > 0$ and $q \in [1/2, 1)$, then the following convergence rates are guaranteed for $\{x^k\}$:*

• If $q = 1/2$, then $\{x^k\}$, $\{\nabla f(x^k)\}$, and $\{f(x^k)\}$ converge linearly to \bar{x} , 0, and $f(\bar{x})$, respectively.

• The setting of $q \in (1/2, 1)$ ensures the estimates

$$\|x^k - \bar{x}\| = \mathcal{O}\left(k^{-\frac{1-q}{2q-1}}\right), \quad \|\nabla f(x^k)\| = \mathcal{O}\left(k^{-\frac{1-q}{2q-1}}\right), \quad \text{and} \quad f(x^k) - f(\bar{x}) = \mathcal{O}\left(k^{-\frac{2-2q}{2q-1}}\right).$$

Proof Since \mathcal{G} is a global approximation of ∇f under condition (3.1), there exists $C > 0$ such that

$$\|\mathcal{G}(x, \delta) - \nabla f(x)\| \leq C\delta \quad \text{for all } (x, \delta) \in \mathbb{R}^n \times (0, \infty). \quad (4.5)$$

By $f \in \mathcal{C}_L^{1,1}$, we find $L > 0$ such that the gradient mapping ∇f is Lipschitz continuous with the constant L on \mathbb{R}^n . Taking the number $N \in \mathbb{N}$ from Proposition 4.2 ensures that $C_k = C_N$ for all $k \geq N$. This implies by Step 2 of Algorithm 1 that

$$f(x^{k+1}) = f\left(x^k - \frac{\kappa}{C_N} g^k\right) \leq f(x^k) - \frac{\kappa(\mu - 2)}{2C_N\mu} \|g^k\|^2 \quad \text{for all } k \geq N, \quad (4.6)$$

which tells us that $\{f(x^k)\}_{k \geq N}$ is decreasing. If $f(x^k) \rightarrow -\infty$, there is nothing to prove, so we assume that $f(x^k) \nrightarrow -\infty$, which implies that $\{f(x^k)\}$ is convergent. As a consequence, we get $f(x^k) - f(x^{k+1}) \rightarrow 0$ as $k \rightarrow \infty$. Then (4.6) tells us that $g^k \rightarrow 0$. From Step 1 of Algorithm 1 it follows that

$$\|g^k\| > \mu C_k \delta_{k+1} = \mu C_N \delta_{k+1} \quad \text{for all } k \geq N \quad (4.7)$$

ensuring that $\delta_{k+1} \downarrow 0$ as $k \rightarrow \infty$. It further follows from $g^k = \mathcal{G}(x^k, \delta_{k+1})$ and (4.5) that

$$\|g^k - \nabla f(x^k)\| = \|\mathcal{G}(x^k, \delta_{k+1}) - \nabla f(x^k)\| \leq C\delta_{k+1} \quad \text{for all } k \in \mathbb{N}, \quad (4.8)$$

which yields $\nabla f(x^k) \rightarrow 0$ as $k \rightarrow \infty$ and thus justifies (i).

To verify (ii), take any accumulation point \bar{x} of $\{x^k\}$ and assume that f satisfies the KL property at \bar{x} . By (4.7) and (4.8), we obtain that

$$\begin{aligned} \|\nabla f(x^k)\| &\leq \|g^k\| + \|\nabla f(x^k) - g^k\| \leq \|g^k\| + C\delta_{k+1} \\ &\leq \|g^k\| + \frac{C\|g^k\|}{\mu C_N} = \alpha \|g^k\| \quad \text{for all } k \geq N, \end{aligned}$$

where $\alpha := \frac{\mu C_N + C}{\mu C_N}$. This together with (4.6) brings us to condition (2.6). By Remark 2.8(i), assumptions (H1) and (H2) in Proposition 2.6 hold. Therefore, $x^k \rightarrow \bar{x}$ as $k \rightarrow \infty$, which justifies (ii).

To proceed with the proof of assertion (iii) under the KL property at \bar{x} with $\psi(t) = Mt^q$, we use the iterations $x^{k+1} = x^k - C^{-1}\kappa g^k$ as in Step 2 of Algorithm 1 together with $\|g^k\| > 0$ from Step 1 of Algorithm 1. This gives us $x^{k+1} \neq x^k$

for $k \geq N$. Combining the latter with (4.6) and (4.9), we see that all the assumptions in Proposition 2.7 are satisfied. This verifies the convergence rates of $\{x^k\}$ to \bar{x} stated in (iii). Since \bar{x} is an accumulation point of $\{x^k\}$, it follows from (i) that \bar{x} is a stationary point of f , i.e., $\nabla f(\bar{x}) = 0$. Hence the usage of Lemma 2.1 and the decreasing property of $\{f(x^k)\}_{k \geq N}$ yields

$$0 \leq f(x^k) - f(\bar{x}) \leq \langle \nabla f(\bar{x}), x^k - \bar{x} \rangle + \frac{L}{2} \|x^k - \bar{x}\|^2 = \frac{L}{2} \|x^k - \bar{x}\|^2,$$

which justifies the convergence rates of $\{f(x^k)\}$ to $f(\bar{x})$ as asserted in (iii).

It remains to verify the convergence rates for $\{\nabla f(x^k)\}$. Since ∇f is Lipschitz continuous with the constant $L > 0$, the claimed property follows from the convergence rates for $\{x^k\}$ due to

$$\|\nabla f(x^k)\| = \|\nabla f(x^k) - \nabla f(\bar{x})\| \leq L \|x^k - \bar{x}\|.$$

This therefore completes the proof of the theorem. \square

Remark 4.4 Here we present a comparison between our analysis for DFC with the analysis in [14]. While both [14] and our paper address the noiseless case and [14] considers a more general approach, our analysis provides additional developments that are not studied in [14]. Specifically:

- (i) Our DFC method (Algorithm 1) explicitly specifies how to construct the gradient approximation. In contrast, [14, Algorithm 3.1] assumes a more general construction and requires the gradient approximation to be sufficiently accurate (as per [14, Assumption 3.1]). In order to make the gradient approximation in DFC satisfying [14, Assumption 3.1], the constant C_k should be larger than C , which is not required in and not ensured by our analysis. Furthermore, the finite difference interval and the stepsize are interacting with each other in our DFC method, while they are considered separately in [14, Algorithm 3.1].
- (ii) Additionally, [14, Algorithm 3.1] employs a different rule for choosing stepsize, allowing it to increase after each iteration, while our DFC does not allow this. The numerical experiments in Subsection 6.1.1 demonstrate that this small change significantly affects the numerical performance of the methods with a more favorable result for DFC.
- (iii) Apart from the differences in algorithmic constructions, our analysis takes a distinct direction by demonstrating the convergence of the gradient sequence to 0 and the convergence of the sequence of iterates to a stationary point. On the other hand, [14, Theorem 3.1] reveals the number of iterations required to reach a near-stationary point, and it also establishes that $\liminf_{k \rightarrow \infty} \|\nabla f(x^k)\| = 0$.

4.3 Analysis for noisy functions

In this part, we provide the convergence analysis with error bounds for DFC in Algorithm 1 addressing problem (1.1) when only a *noisy approximation* $\phi(x) =$

$f(x) + \xi(x)$ of f is available, where $\xi : \mathbb{R}^n \rightarrow \mathbb{R}$ is a *noise function* bounded by some constant $\xi_f > 0$, i.e.,

$$|\xi(x)| \leq \xi_f \text{ for all } x \in \mathbb{R}^n. \quad (4.9)$$

Due to the design of DFC, we *do not* assume that ξ_f is known. For brevity, consider only the *forward finite difference approximation*, while other gradient approximation methods can be employed via modifications of the inexact conditions in Definition 3.1 for noisy functions. We first construct the gradient approximation for f via the forward finite difference with the noisy function ϕ defined by

$$\tilde{\mathcal{G}}(x, \delta) := \frac{1}{\delta} \sum_{i=1}^n (\phi(x + \delta e_i) - \phi(x)) e_i \text{ for any } (x, \delta) \in \mathbb{R}^n \times (0, \infty). \quad (4.10)$$

Recall the following noisy version of Proposition 3.4, which is well known and can be found in, e.g., [7, Theorem 2.1].

Proposition 4.5 *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a $C_L^{1,1}$ function. Then the noisy forward finite difference (4.10) satisfies the error bound*

$$\|\tilde{\mathcal{G}}(x, \delta) - \nabla f(x)\| \leq \frac{L\sqrt{n}\delta}{2} + \frac{2\sqrt{n}\xi_f}{\delta} \text{ for all } \delta > 0. \quad (4.11)$$

For a better exposition, consider DFC with specific parameters $\mu = 4$ and $\kappa = \frac{\sqrt{n}}{2}$, although other general selections of $\mu > 2$ and $\kappa > 0$ still work with the same analysis. We also define $L_k := \frac{C_k}{\kappa}$ for each $k \in \mathbb{N}$ as approximate Lipschitz constants. In order to deal with noise, a relaxation is required in the descent condition in Step 2 of DFC, which leads us to the following algorithm.

Algorithm 2 (DFC for noisy functions).

Step 0 (initialization). Select some $x^1 \in \mathbb{R}^n$, $\delta_1 > 0$, $L_1 > 0$, $\theta \in (0, 1)$, and $\eta > 1$.

Step 1 (approximate gradient). Find g^k and the smallest nonnegative integer i_k such that

$$g^k = \tilde{\mathcal{G}}(x^k, \theta^{i_k} \delta_k) \text{ and } \|g^k\| > 2L_k \sqrt{n} \theta^{i_k} \delta_k. \quad (4.12)$$

Then set $\delta_{k+1} := \theta^{i_k} \delta_k$.

Step 2 (update). If $\phi\left(x^k - \frac{1}{L_k} g^k\right) \leq \phi(x^k) - \frac{1}{24L_k} \|g^k\|^2$, then $x^{k+1} := x^k - \frac{1}{L_k} g^k$ and $L_{k+1} := L_k$. Otherwise, $x^{k+1} := x^k$ and $L_{k+1} := \eta L_k$.

In the following remark, we present general ideas for the construction of Algorithm 2 and discuss the main differences between DFC and derivative-free trust-region methods presented in [17, Section 10.3].

Remark 4.6 (i) (*General ideas*) The ideas for the construction of DFC are as follows.

Firstly, since the noise is only bounded and may not vanish, we may not make any improvement if the k^{th} iterate x^k is sufficiently close to the optimal solutions. However, if x^k is far from the optimal solutions, we can expect its function value $f(x^k)$ to be sufficiently larger than the noise. In this case, the total error for approximating the gradient is dominated by the truncated one, i.e.,

$$\|g^k - \nabla f(x^k)\| \leq \frac{L\sqrt{n}\delta}{2} + \frac{2\sqrt{n}\xi_f}{\delta} \approx \frac{L\sqrt{n}\delta}{2} \text{ for } \delta > 0. \quad (4.13)$$

Next we try to find the largest finite difference interval δ_{k+1} such that the corresponding approximation is acceptable in the sense that

$$\|g^k - \nabla f(x^k)\| \leq \frac{1}{4} \|g^k\|, \quad g^k = \tilde{\mathcal{G}}(x^k, \delta_{k+1}),$$

which ensures a sufficient decrease of the function f when moving towards the direction $-g^k$. Note that another condition that makes the \approx sign in (4.13) reliable is that δ should be sufficiently large; otherwise, the roundoff error becomes dominant again. This is where the smallest property of i_k plays an important role and motivates us to consider construction (4.12) in Step 1. The reasons behind using L_k to approximate L in Step 2 are exactly the same as using C_k to approximate C discussed in Remark 4.1 for the noiseless case.

- (ii) (*Comparison with trust-region algorithms*) In contrast to the derivative-free trust-region algorithms presented in [17, Section 10.3], we do not construct approximation models. Instead, we directly use the finite-difference gradient as the descent direction and move along the latter. Additionally, the accuracy of the approximate gradient and the stepsize in DFC are updated in separate steps, in contrast to the fact that these updates are combined into the trust-region radius in the methods described in [17, Section 10.3].

Now we start deriving fundamental properties of Algorithm 2 that are essential for the convergence analysis. Due to the presence of noise, there is no guarantee that Step 1 in Algorithm 2 will terminate after a finite number of trials for i_k . Therefore, we say that Step 1 is *successful* if such i_k is found, and is *unsuccessful* otherwise.

Note that the minimum of the right-hand side of (4.11) is achieved when $\delta = \sqrt{\frac{4\xi_f}{L}}$, which we regard as the (unknown) optimal finite difference interval. Given this, we begin our analysis with a result that guarantees the success of Step 1 in Algorithm 2 provided that the gradient of the current iterate is not close to 0, that the finite difference interval above is the optimal threshold, and that the approximate Lipschitz constant is sufficiently small.

Proposition 4.7 *At the k^{th} iteration of Algorithm 2, if the conditions*

$$\|\nabla f(x^k)\| \geq (4\theta^{-1}\eta + \theta^{-1} + 1)\sqrt{Ln\xi_f}, \quad \delta_k \geq \sqrt{\frac{4\xi_f}{L}}, \quad L_k < \eta L$$

are satisfied, then Step 1 is successful with $\delta_{k+1} \geq \sqrt{\frac{4\xi_f}{L}}$.

Proof Let $i := \lfloor \log_{\theta} \left(\frac{1}{\delta_k} \sqrt{\frac{4\xi_f}{L}} \right) \rfloor$, where $\lfloor \cdot \rfloor$ stands for the floor/greatest integer function. Since we have $\delta_k \geq \sqrt{\frac{4\xi_f}{L}}$ and $\theta \in (0, 1)$, the number i is a nonnegative integer satisfying the inclusion

$$i \in \left(\log_{\theta} \left(\frac{1}{\delta_k} \sqrt{\frac{4\xi_f}{L}} \right) - 1, \log_{\theta} \left(\frac{1}{\delta_k} \sqrt{\frac{4\xi_f}{L}} \right) \right]$$

while implying in turn that $\theta^i \delta_k \in \left[\sqrt{\frac{4\xi_f}{L}}, \theta^{-1} \sqrt{\frac{4\xi_f}{L}} \right)$. We now show that inequality (4.12) in Step 1 of Algorithm 2 is satisfied for $i_k = i$, which yields the success of the step with $\delta_{k+1} \geq \theta^i \delta_k \geq \sqrt{\frac{4\xi_f}{L}}$. Indeed, with $g_i^k := \tilde{\mathcal{G}}(x^k, \theta^i \delta_k)$, the error bound (4.11) and $\theta^i \delta_k \in \left[\sqrt{\frac{4\xi_f}{L}}, \theta^{-1} \sqrt{\frac{4\xi_f}{L}} \right)$ tell us that

$$\begin{aligned} \|g_i^k - \nabla f(x^k)\| &\leq \frac{L\sqrt{n}}{2} \theta^i \delta_k + \frac{2\sqrt{n}\xi_f}{\theta^i \delta_k} \\ &\leq \theta^{-1} \sqrt{Ln\xi_f} + \sqrt{Ln\xi_f} = (1 + \theta^{-1})\sqrt{Ln\xi_f}. \end{aligned}$$

Since $\|\nabla f(x^k)\| \geq (4\theta^{-1}\eta + \theta^{-1} + 1)\sqrt{Ln\xi_f}$ and $L_k < \eta L$, we get that

$$\begin{aligned} \|g_i^k\| &\geq \|\nabla f(x^k)\| - \|g_i^k - \nabla f(x^k)\| \geq 4\theta^{-1}\eta\sqrt{Ln\xi_f} = 2\eta L\sqrt{n}\theta^{-1}\sqrt{\frac{4\xi_f}{L}} \\ &> 2L_k\sqrt{n}\theta^i \delta_k. \end{aligned}$$

Therefore, Step 1 of Algorithm 2 is successful with $\delta_{k+1} \geq \theta^i \delta_k \geq \sqrt{\frac{4\xi_f}{L}}$. \square

Proposition 4.8 *At the k^{th} iteration of Algorithm 2, suppose that $\|\nabla f(x^k)\| \geq 16\sqrt{Ln\xi_f}$ and that Step 1 is successful with $\delta_{k+1} \geq \sqrt{\frac{4\xi_f}{L}}$. The following assertions hold:*

- (i) $\|g^k - \nabla f(x^k)\| \leq \frac{4L+L_k}{15L_k} \|g^k\|$.
- (ii) If in addition $L_k \geq L$, then $\langle \nabla f(x^k), g^k \rangle \geq \frac{2}{3} \|g^k\|^2$ and $\|\nabla f(x^k)\| \leq \frac{4}{3} \|g^k\|$.
- (iii) If in addition $L_k < \eta L$, then $L_{k+1} < \eta L$.

Proof (i) Since Step 1 of Algorithm 2 is successful, we deduce that $\|g^k\| \geq 2L_k\sqrt{n}\delta_{k+1}$. Combining this with $g^k = \tilde{G}(x^k, \delta_{k+1})$, (4.11), and the estimates $\delta_{k+1} \geq \sqrt{\frac{4\xi_f}{L}}$, $\|\nabla f(x^k)\| \geq 16\sqrt{Ln\xi_f}$ brings us to

$$\begin{aligned} \|g^k - \nabla f(x^k)\| &\leq \frac{L\sqrt{n}}{2}\delta_{k+1} + \frac{2\sqrt{n}\xi_f}{\delta_{k+1}} \\ &\leq \frac{L\sqrt{n}}{2} \frac{\|g^k\|}{2L_k\sqrt{n}} + \sqrt{Ln\xi_f} \\ &\leq \frac{L}{4L_k} \|g^k\| + \frac{1}{16} \|\nabla f(x^k)\| \\ &\leq \frac{L}{4L_k} \|g^k\| + \frac{1}{16} \|g^k\| + \frac{1}{16} \|g^k - \nabla f(x^k)\|. \end{aligned}$$

The latter inequality yields $\|g^k - \nabla f(x^k)\| \leq \frac{4L+L_k}{15L_k} \|g^k\|$, which verifies (i).

(ii) Using (i) and $L_k \geq L$ gives us $\|g^k - \nabla f(x^k)\| \leq \frac{1}{3} \|g^k\|$. Combining this estimates with the Cauchy-Schwarz inequality, we arrive at

$$\begin{aligned} \langle \nabla f(x^k), g^k \rangle &\geq \langle \nabla f(x^k) - g^k, g^k \rangle + \|g^k\|^2 \\ &\geq -\|\nabla f(x^k) - g^k\| \|g^k\| + \|g^k\|^2 \geq \frac{2}{3} \|g^k\|^2. \end{aligned}$$

In addition, it follows from $\|g^k - \nabla f(x^k)\| \leq \frac{1}{3} \|g^k\|$ that

$$\|g^k\| \geq \|\nabla f(x^k)\| - \|g^k - \nabla f(x^k)\| \geq \|\nabla f(x^k)\| - \frac{1}{3} \|g^k\|,$$

which justifies the claimed estimates $\|\nabla f(x^k)\| \leq \frac{4}{3} \|g^k\|$.

(iii) This assertion obviously holds if $L_{k+1} = L_k$, so we consider the case where $L_{k+1} = \eta L_k$. Suppose on the contrary that $L_{k+1} \geq \eta L$, which yields $L_k \geq L$. Then it follows from assertion (ii) that $\langle \nabla f(x^k), g^k \rangle \geq \frac{2}{3} \|g^k\|^2$. Furthermore, Lemma 2.1 tells us that

$$\begin{aligned} f\left(x^k - \frac{1}{L_k} g^k\right) - f(x^k) &\leq -\frac{1}{L_k} \langle \nabla f(x^k), g^k \rangle + \frac{L}{2} \left\| \frac{1}{L_k} g^k \right\|^2 \\ &\leq -\frac{1}{L_k} \frac{2}{3} \|g^k\|^2 + \frac{1}{2L_k} \|g^k\|^2 = -\frac{1}{6L_k} \|g^k\|^2. \end{aligned} \quad (4.14)$$

Since Step 1 is successful with $\delta_{k+1} \geq \sqrt{\frac{4\xi_f}{L}}$ and $L_k \geq L$, we deduce that

$$\|g^k\| > 2L_k\delta_{k+1} \geq 2L_k\sqrt{\frac{4\xi_f}{L}} \geq 4\sqrt{L_k\xi_f},$$

which means that $\xi_f \leq \frac{1}{16L_k} \|g^k\|^2$. Combining the latter with (4.9) and (4.14) gives us

$$\begin{aligned} \phi\left(x^k - \frac{1}{L_k}g^k\right) - \phi(x^k) &\leq f\left(x^k - \frac{1}{L_k}g^k\right) - f(x^k) + 2\xi_f \\ &\leq -\frac{1}{6L_k} \|g^k\|^2 + \frac{1}{8L_k} \|g^k\|^2 = -\frac{1}{24L_k} \|g^k\|^2. \end{aligned}$$

By Step 2 of Algorithm 2, it follows that $L_{k+1} = L_k$, a contradiction, which justifies $L_{k+1} < \eta L$. \square

In the propositions above, we can choose $\theta \in (0, 1)$ and $\eta > 1$ to get $4\theta^{-1}\eta + \theta^{-1} + 1 < 16$ by taking, e.g., $\theta = \frac{\sqrt{2}}{2}$ and $\eta = 2$. To simplify the presentation, we make such a *selection of parameters* in the results below. Under this choice, the following property of Algorithm 2 can be deduced immediately from Proposition 4.7 and Proposition 4.8(iii).

Proposition 4.9 *Let $\{x^k\}$ be generated by Algorithm 2 with $L_1 < \eta L$ and $\delta_1 \geq \sqrt{\frac{4\xi_f}{L}}$. If for some $K \in \mathbb{N}$ we have $\|\nabla f(x^k)\| \geq 16\sqrt{Ln\eta\xi_f}$ whenever $k = 1, \dots, K$, then Step 1 is successful with $\delta_{k+1} \geq \sqrt{\frac{4\xi_f}{L}}$ and $L_{k+1} < \eta L$ for all $k = 1, \dots, K$.*

Now we are ready to establish the main convergence properties of Algorithm 2.

Theorem 4.10 *Let $\{x^k\}$ be generated by Algorithm 2 with $\delta_1 \geq \sqrt{\frac{4\xi_f}{L}}$ and $L_1 < \eta L$. Then the number N of iterations that Algorithm 2 takes until $\|\nabla f(x^N)\| < 16\sqrt{Ln\xi_f}$ is bounded by*

$$\begin{aligned} N \leq N_{\text{opt}} &:= 1 + \left\lfloor \frac{f(x^1) - f^* + 2\xi_f}{M\xi_f} \right\rfloor + \left\lfloor \log_\eta \left(\frac{\eta L}{L_1} \right) \right\rfloor, \text{ where} \\ M &:= \frac{150nL_1^2}{\eta(L + 4L_1)^2} \text{ and } f^* := \inf_{x \in \mathbb{R}^n} f(x) > -\infty. \end{aligned}$$

The total number N_{fval} of function evaluations needed to achieve this goal is bounded by

$$N_{\text{fval}} \leq (n + 2)N_{\text{opt}} + n \left\lfloor \log_\theta \left(\frac{2\sqrt{\xi_f}}{\delta_1\sqrt{L}} \right) \right\rfloor.$$

Proof If Step 1 is unsuccessful for the first time at $K \leq N_{\text{opt}}$, then it follows from Proposition 4.9 that $\|\nabla f(x^N)\| < 16\sqrt{Ln\xi_f}$ for some $N \leq K$, which verifies the claimed bound. Now we suppose that Step 1 is successful for all $k = 1, \dots, N_{\text{opt}}$ and assume on the contrary that

$$\|\nabla f(x^k)\| \geq 16\sqrt{Ln\xi_f} \text{ for all } k = 1, \dots, N_{\text{opt}}.$$

Proposition 4.9 tells us that $\delta_{k+1} \geq \sqrt{\frac{4\xi_f}{L}}$ and $L_{k+1} < \eta L$ for all $k = 1, \dots, N_{\text{opt}}$, and thus it follows from Proposition 4.8(i) and the construction of $\{L_k\}$ that

$$\|g^k - \nabla f(x^k)\| \leq \frac{4L + L_k}{15L_k} \|g^k\| \leq \frac{4L + L_1}{15L_1} \|g^k\|.$$

This gives us in turn the estimates

$$16\sqrt{Ln\xi_f} \leq \|\nabla f(x^k)\| \leq \|g^k\| + \|g^k - \nabla f(x^k)\| \leq \frac{4L + 16L_1}{15L_1} \|g^k\| \text{ for all } k = 1, \dots, N_{\text{opt}}. \quad (4.15)$$

Define $I := \{k \in \mathbb{N} \mid 1 \leq k \leq N_{\text{opt}}, L_{k+1} = L_k\}$ and deduce from the construction of $\{L_k\}$ with $L_{k+1} < \eta L$ as $k = 1, \dots, N_{\text{opt}}$ that there are at most $\left\lfloor \log_\eta \left(\frac{\eta L}{L_1} \right) \right\rfloor$ iterations for which $L_{k+1} = \eta L_k$. This yields

$$|I| \geq N_{\text{opt}} - \left\lfloor \log_\eta \left(\frac{\eta L}{L_1} \right) \right\rfloor = 1 + \left\lfloor \frac{f(x^1) - f^* + 2\xi_f}{M\xi_f} \right\rfloor. \quad (4.16)$$

Take any $k = 1, \dots, N_{\text{opt}}$. If $k \notin I$, we get that $\phi(x^{k+1}) = \phi(x^k)$. For any $k \in I$, it follows from Step 2 of Algorithm 2, $L_k < \eta L$ and (4.15) that

$$\phi(x^{k+1}) - \phi(x^k) \leq -\frac{1}{24L_k} \|g^k\|^2 \leq -\frac{1}{24\eta L} \|g^k\|^2 \leq -M\xi_f,$$

where M is defined in the statement of the theorem. Since $\phi(x^k) = \phi(x^{k+1})$ when $k \notin I$, we have

$$f^* - \xi_f \leq \phi(x^{N_{\text{opt}}+1}) = \phi(x^1) + \sum_{k=1}^{N_{\text{opt}}} (\phi(x^{k+1}) - \phi(x^k)) \leq f(x^1) + \xi_f - |I|M\xi_f,$$

which yields $|I| \leq \frac{f(x^1) - f^* + 2\xi_f}{M\xi_f}$ and thus contradicts (4.16).

(iii) By (ii), at most N_{opt} iterations are needed to reach the near stationary point. For each iteration of Algorithm 2, we need at least one approximate gradient evaluation g^k in Step 1. Since $\{\delta_k\}$ is nonincreasing with $\delta_1 \geq \delta_k \geq \sqrt{\frac{4\xi_f}{L}}$, the number of additional approximate gradient evaluations g^k required to adjust the finite difference intervals δ_k throughout all the iterations is at most $\left\lceil \log_{\theta} \left(\frac{2\sqrt{\xi_f}}{\delta_1\sqrt{L}} \right) \right\rceil$. Employing the forward finite difference, we can reuse $\phi(x^k)$ for additional gradient evaluations. This tells us that the total number of function evaluations for determining the approximate gradient g^k is at most $(n+1)N_{\text{opt}} + n \left\lceil \log_{\theta} \left(\frac{2\sqrt{\xi_f}}{\delta_1\sqrt{L}} \right) \right\rceil$. We also need one additional function evaluation to check the descent condition in Step 2 of Algorithm 2 at each iteration, which results in at most $(n+2)N_{\text{opt}} + n \left\lceil \log_{\theta} \left(\frac{2\sqrt{\xi_f}}{\delta_1\sqrt{L}} \right) \right\rceil$ total function evaluations. \square

Remark 4.11 Let us briefly discuss relationships between our analysis for DFC and the analysis in [6]. First observe that the noise level is unknown for our DFC algorithm, while it is required to be known for the analysis in [6] as mentioned after [6, Assumption 1.3]. The algorithmic construction of [6, Algorithm 2.1] shares many similarities with [14, Algorithm 3.1], and so it has some major differences with our DFC as mentioned above in Remark 4.4(i,ii).

5 General derivative-free methods for $C^{1,1}$ functions

In this section, we consider problem (1.1), where f is of class $C^{1,1}$ and develop new derivative-free optimization methods in both cases of noiseless and noisy objective functions.

5.1 Backtracking linesearch for noiseless functions

Here we propose and justify the novel *derivative-free method with backtracking step-size* (DFB) to solve the optimization problem (1.1) in the noiseless setting. The main result of this subsection establishes the *global convergence* with convergence rates of the following algorithm, which employs gradient approximations satisfying the local error bound estimate (3.2).

Algorithm 3 (DFB).

Step 0 (initialization). Choose a local approximation \mathcal{G} of ∇f under condition (3.2). Select an initial point $x^1 \in \mathbb{R}^n$ and initial radius $\delta_1 > 0$, a constant $C_1 > 0$, factors $\theta \in (0, 1)$, $\mu > 2$, $\eta > 1$, linesearch constants $\beta \in (0, 1/2)$, $\gamma \in (0, 1)$, $\bar{\tau} > 0$, and an initial bound $t_1^{\min} \in (0, \bar{\tau})$. Choose a sequence of manually controlled errors $\{v_k\} \subset [0, \infty)$ such that $v_k \downarrow 0$ as $k \rightarrow \infty$. Set $k := 1$.

Step 1 (approximate gradient). Select g^k and the smallest nonnegative integer i_k so that

$$g^k = \mathcal{G}(x^k, \min\{\theta^{i_k} \delta_k, v_k\}) \quad \text{and} \quad \|g^k\| > \mu C_k \theta^{i_k} \delta_k. \quad (5.1)$$

Then set $\delta_{k+1} := \theta^{i_k} \delta_k$.

Step 2 (linesearch). Set the tentative stepsize $t_k := \bar{\tau}$. While

$$f(x^k - t_k g^k) > f(x^k) - \beta t_k \|g^k\|^2 \quad \text{and} \quad t_k \geq t_k^{\min}, \quad (5.2)$$

set $t_k := \gamma t_k$.

Step 3 (stepsize and parameters update). If $t_k \geq t_k^{\min}$, then set $\tau_k := t_k$, $C_{k+1} := C_k$, and $t_{k+1}^{\min} := t_k^{\min}$. Otherwise, set $\tau_k := 0$, $C_{k+1} := \eta C_k$, and $t_{k+1}^{\min} := \gamma t_k^{\min}$.

Step 4 (iteration update). Set $x^{k+1} := x^k - \tau_k g^k$. Increase k by 1 and go back to Step 1.

Remark 5.1 (i) Fix any $k \in \mathbb{N}$. The procedure of finding g^k and i_k that satisfies Step 1 of Algorithm 3 can be described as follows. Set $i_k := 0$ and calculate g^k as

$$g^k = \mathcal{G}(x^k, \min\{\theta^{i_k} \delta_k, v_k\}). \quad (5.3)$$

While $\|g^k\| \leq \mu C_k \theta^{i_k} \delta_k$, increase i_k by 1 and recalculate g^k by formula (5.3). We now show that when $\nabla f(x^k) \neq 0$, this procedure stops after a *finite number of steps* giving us g^k and i_k as desired. Indeed, assuming on the contrary that the procedure does not stop, we get a sequence of $\{g_i^k\}$ with

$$g_i^k = \mathcal{G}(x^k, \min\{\theta^i \delta_k, v_k\}) \quad \text{and} \quad \|g_i^k\| \leq \mu C_k \theta^i \delta_k \quad \text{for all } i \in \mathbb{N}. \quad (5.4)$$

Since \mathcal{G} is a local approximation of ∇f , for any fixed $\Delta > 0$ condition (3.2) with $\Omega = \{x^k\}$ ensures the existence of a positive number C such that

$$\|\mathcal{G}(x^k, \delta) - \nabla f(x^k)\| \leq C\delta \quad \text{whenever } 0 < \delta \leq \Delta. \quad (5.5)$$

By $\theta \in (0, 1)$, there is $N \in \mathbb{N}$ with $\theta^i \delta_k \leq \Delta$ for all $i \geq N$. Combining this with (5.4) and (5.5) yields

$$\|g_i^k - \nabla f(x^k)\| \leq C\theta^i \delta_k \quad \text{and} \quad \|g_i^k\| \leq \mu C_k \theta^i \delta_k \quad \text{for all } i \geq N.$$

Letting $i \rightarrow \infty$, we arrive at $\nabla f(x^k) = 0$, which is a contradiction.

(ii) It follows directly from the construction of δ_k in Step 1 of Algorithm 3 that

$$g^k = \mathcal{G}(x^k, \min\{\delta_{k+1}, v_k\}) \quad \text{and} \quad \|g^k\| > \mu C_k \delta_{k+1}. \quad (5.6)$$

To proceed further with the convergence analysis of Algorithm 3, we obtain two results of their independent interest. The first one reveals some *uniformity* of general linesearch procedures with respect to the selections of reference points, stepsizes, and directions.

Lemma 5.2 *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a function with a locally Lipschitz continuous gradient, and let $\beta \in (0, 1/2)$. Then for any nonempty bounded set $\Omega \subset \mathbb{R}^n$, there exists $\bar{t} > 0$ such that*

$$f(x - tg) \leq f(x) - \beta t \|g\|^2 \quad \text{whenever } x \in \Omega, \quad 2\|g - \nabla f(x)\| \leq \|g\|, \\ \text{and } t \in (0, \bar{t}].$$

Proof The boundedness of Ω gives us $r > 0$ such that $\Omega \subset r\overline{\mathbb{B}}$. Using the continuity of ∇f and the compactness of $r\overline{\mathbb{B}}$, define $r' := \max\{\|\nabla f(x)\| \mid x \in r\overline{\mathbb{B}}\}$. Since $f \in C^{1,1}$, there exists $L > 0$ such that ∇f is Lipschitz continuous with the constant L on $\Theta := (r + 2r')\overline{\mathbb{B}}$. By $\beta < 1/2$, we find $\bar{t} > 0$ with $\bar{t} < \min\{1, L^{-1}(1 - 2\beta)\}$. Now take some $x \in \Omega \subset \Theta$ and $g \in \mathbb{R}^n$ such that $2\|g - \nabla f(x)\| \leq \|g\|$ and $t \in (0, \bar{t}]$. The choice of g gives us by the Cauchy-Schwarz inequality that

$$\langle \nabla f(x), g \rangle = \langle \nabla f(x) - g, g \rangle + \|g\|^2 \geq -\|\nabla f(x) - g\| \|g\| + \|g\|^2 \\ \geq -\frac{1}{2} \|g\|^2 + \|g\|^2 = \frac{1}{2} \|g\|^2, \quad (5.7)$$

and by using the triangle inequality that

$$\|\nabla f(x)\| \geq \|g\| - \|g - \nabla f(x)\| \geq \frac{1}{2} \|g\|.$$

Combining the latter with the choice of t , \bar{t} , $x \in \Omega \subset r\overline{\mathbb{B}}$ and the construction of r' yields

$$t \|g\| \leq \bar{t} \|g\| \leq 2\bar{t} \|\nabla f(x)\| \leq 2\bar{t}r' < 2r',$$

which ensures that $x - tg \in \Theta$. The convexity of Θ tells us that the entire line segment $[x, x - tg]$ lies on Θ . Remembering that ∇f is Lipschitz continuous with the constant

L on Θ , we employ Lemma 2.1 by taking into account that $t \leq \bar{t} < L^{-1}(1 - 2\beta)$ and that (5.7). This gives us

$$\begin{aligned} f(x - tg) - f(x) &\leq \langle x - tg - x, \nabla f(x) \rangle + \frac{L}{2} \|x - tg - x\|^2 \\ &= -t \langle g, \nabla f(x) \rangle + \frac{Lt^2}{2} \|g\|^2 \leq -\frac{t}{2} \|g\|^2 + \frac{Lt^2}{2} \|g\|^2 \\ &= -\beta t \|g\|^2 + t \|g\|^2 \frac{2\beta - 1 + Lt}{2} \leq -\beta t \|g\|^2 \end{aligned}$$

and thus completes the proof of the lemma. \square

Employing the obtained lemma, we derive the next result showing that unless the stationary point is found, Algorithm 3 always makes a progress after a finite number of iterations.

Proposition 5.3 *Let $\{x^k\}$ and $\{\tau_k\}$ be the sequences generated by Algorithm 3, and let $K \in \mathbb{N}$ be such that $\nabla f(x^K) \neq 0$. Then we can choose a number $N \geq K$ so that $\tau_N > 0$.*

Proof Assume on the contrary that $\tau_k = 0$ for all $k \geq K$. Steps 3 and 4 of Algorithm 3 give us

$$t_{k+1}^{\min} = \gamma t_k^{\min} \quad \text{and} \quad x^k = x^K \quad \text{for all } k \geq K. \quad (5.8)$$

Therefore, $\nabla f(x^k) = \nabla f(x^K) \neq 0$ for all $k \geq K$, which implies that g^k and i_k in Step 1 of Algorithm 3 exist for all $k \geq K$. Since \mathcal{G} is a local approximation of ∇f , for any fixed $\Delta > 0$ condition (3.2) with $\Omega = \{x^K\}$ ensures the existence of $C > 0$ with

$$\|\mathcal{G}(x^K, \delta) - \nabla f(x^K)\| \leq C\delta \quad \text{whenever } 0 < \delta \leq \Delta. \quad (5.9)$$

It follows from Lemma 5.2 with $\Omega = \{x^K\}$ that there exists some $\bar{t} > 0$ such that

$$\begin{aligned} f(x^K - tg) &\leq f(x^K) - \beta t \|g\|^2 \quad \text{whenever } 2 \|g - \nabla f(x^K)\| \leq \|g\| \\ &\text{and } t \in (0, \bar{t}]. \end{aligned} \quad (5.10)$$

Using $v_k \downarrow 0$, $t_k^{\min} \downarrow 0$, $\nabla f(x^K) \neq 0$, and (5.8) gives us $N \geq K$ for which $v_N < \min \left\{ \Delta, \frac{1}{3C} \|\nabla f(x^K)\| \right\}$ and $t_N^{\min} < \gamma \bar{t}$. Then we get from (5.9) with taking into account $x^N = x^K$ that

$$\|\mathcal{G}(x^N, \min \{\delta_{N+1}, v_N\}) - \nabla f(x^N)\| \leq C \min \{\delta_{N+1}, v_N\} \leq C v_N \leq \frac{1}{3} \|\nabla f(x^N)\|.$$

Combining this with $g^N = \mathcal{G}(x^N, \min \{\delta_{N+1}, v_N\})$ from (5.6) provides the estimate

$$\|g^N - \nabla f(x^N)\| \leq \frac{1}{3} \|\nabla f(x^N)\|,$$

which implies by the triangle inequality that

$$\|g^N\| \geq \|\nabla f(x^N)\| - \|g^N - \nabla f(x^N)\| \geq 2\|g^N - \nabla f(x^N)\|.$$

Employing the latter together with (5.10) and $x^N = x^K$ yields

$$f(x^N - tg^N) \leq f(x^N) - \beta t \|g^N\|^2 \text{ for all } t \in (0, \bar{t}]. \quad (5.11)$$

It follows from (5.11) and the choice of parameters that

$$\max \{t \mid f(x^N - tg^N) \leq f(x^N) - \beta t \|g^N\|^2, t = \bar{t}, \bar{t}\gamma, \bar{t}\gamma^2, \dots\} > \gamma\bar{t} > t_N^{\min},$$

which implies in turn by Step 2 of Algorithm 3 that

$$t_N = \max \{t \mid f(x^N - tg^N) \leq f(x^N) - \beta t \|g^N\|^2, t = \bar{t}, \bar{t}\gamma, \bar{t}\gamma^2, \dots\} > t_N^{\min}.$$

By Step 3 of Algorithm 3, we conclude that $\tau_N = t_N > 0$, a contradiction completing the proof. \square

Now we are ready to establish the convergence properties of Algorithm 3.

Theorem 5.4 *Let $\{x^k\}$ be the sequence generated by Algorithm 3 and assume that $\nabla f(x^k) \neq 0$ for all $k \in \mathbb{N}$. Then either $f(x^k) \rightarrow -\infty$ as $k \rightarrow \infty$, or the following assertions hold:*

- (i) *Every accumulation point of $\{x^k\}$ is a stationary point of f .*
- (ii) *If the sequence $\{x^k\}$ is bounded, then the set of accumulation points of $\{x^k\}$ is nonempty, compact, and connected in \mathbb{R}^n .*
- (iii) *If $\{x^k\}$ has an isolated accumulation point \bar{x} , then this sequence converges to \bar{x} .*

Proof First it follows from Steps 2 and 3 of Algorithm 3 that

$$\beta\tau_k \|g^k\|^2 \leq f(x^k) - f(x^{k+1}) \text{ for all } k \in \mathbb{N}, \quad (5.12)$$

which tells us that $\{f(x^k)\}$ is nonincreasing. If $f(x^k) \rightarrow -\infty$, there is nothing to prove; so we assume that $f(x^k) \nrightarrow -\infty$, which implies by the nonincreasing property of $\{f(x^k)\}$ that $\inf f(x^k) > -\infty$. Summing up the inequalities in (5.12) over $k = 1, 2, \dots$ with taking into account that $x^{k+1} = x^k - \tau_k g^k$ from the update in Step 3 of Algorithm 3 gives us

$$\sum_{k=1}^{\infty} \tau_k \|g^k\|^2 < \infty \text{ and } \sum_{k=1}^{\infty} \|g^k\| \cdot \|x^{k+1} - x^k\| < \infty. \quad (5.13)$$

We divide the proof of (i) into two parts by showing first that the origin is an accumulation point of $\{g^k\}$ and then employing Lemma 2.2 to establish stationarity of all the accumulation points of $\{x^k\}$.

Claim 1 *The origin $0 \in \mathbb{R}^n$ is an accumulation point of the sequence $\{g^k\}$.*

Arguing by contradiction, suppose that there are numbers $\varepsilon > 0$ and $K \in \mathbb{N}$ such that

$$\|g^k\| \geq \varepsilon \text{ for all } k \geq K. \quad (5.14)$$

Combining this with (5.13) gives us $\tau_k \downarrow 0$ and $\sum_{k=1}^{\infty} \|x^{k+1} - x^k\| < \infty$. The latter implies that $\{x^k\}$ converges to some $\bar{x} \in \mathbb{R}^n$. By taking a larger K , we can assume that $\tau_k < \bar{\tau}$ for all $k \geq K$. Let \mathcal{N} be the set of all $k \in \mathbb{N}$ such that $\tau_k > 0$. It follows from Proposition 5.3 that \mathcal{N} is infinite. Hence we can take any $k \geq K$ with $k \in \mathcal{N}$ and get that $\tau_k \in (0, \bar{\tau})$. Step 3 of Algorithm 3 ensures that $\tau_k = t_k \in [t_k^{\min}, \bar{\tau})$. Fixing such an index k , we get from the exit condition in Step 2 of Algorithm 3 that

$$-\gamma^{-1}\beta\tau_k \|g^k\|^2 < f(x^k - \gamma^{-1}\tau_k g^k) - f(x^k). \quad (5.15)$$

The classical mean value theorem gives us $\tilde{x}^k \in [x^k, x^k - \gamma^{-1}\tau_k g^k]$ such that

$$f(x^k - \gamma^{-1}\tau_k g^k) - f(x^k) = -\gamma^{-1}\tau_k \langle g^k, \nabla f(\tilde{x}^k) \rangle. \quad (5.16)$$

Combining this with (5.15) yields

$$-\gamma^{-1}\beta\tau_k \|g^k\|^2 < -\gamma^{-1}\tau_k \langle g^k, \nabla f(\tilde{x}^k) \rangle,$$

which implies by dividing both sides of the inequality by $-\gamma^{-1}\tau_k < 0$ that

$$\langle g^k, \nabla f(\tilde{x}^k) \rangle < \beta \|g^k\|^2 \text{ for all } k \geq K, k \in \mathcal{N}. \quad (5.17)$$

Take some neighborhood Ω of \bar{x} and $\Delta > 0$. Since \mathcal{G} is a local approximation of ∇f under condition (3.2), there exists $C > 0$ such that

$$\|\mathcal{G}(x, \delta) - \nabla f(x)\| \leq C\delta \text{ whenever } 0 < \delta \leq \Delta \text{ and } x \in \Omega. \quad (5.18)$$

Since $v_k \downarrow 0$ and $x^k \rightarrow \bar{x}$, by taking a larger K we can assume that $v_k < \Delta$ and $x^k \in \Omega$ for all $k \geq K$. Using this together with (5.18) and $g^k = \mathcal{G}(x^k, \min\{\delta_{k+1}, v_k\})$ in (5.6) tells us that

$$\|g^k - \nabla f(x^k)\| = \|\mathcal{G}(x^k, \min\{\delta_{k+1}, v_k\}) - \nabla f(x^k)\| \leq C \min\{\delta_{k+1}, v_k\} \leq C v_k.$$

Combining the latter with $x^k \rightarrow \bar{x}$, $v_k \downarrow 0$ as $k \rightarrow \infty$, and the continuity of ∇f gives us

$$g^k \rightarrow \nabla f(\bar{x}) \text{ as } k \rightarrow \infty, \quad (5.19)$$

which yields $\|\nabla f(\bar{x})\| > 0$ by (5.14). It follows from (5.19), $\tau_k \downarrow 0$, $x^k \rightarrow \bar{x}$, and $\tilde{x}^k \in [x^k, x^k - \gamma^{-1}\tau_k g^k]$ for all $k \geq K$ with $k \in \mathcal{N}$ that $\tilde{x}^k \xrightarrow{\mathcal{N}} \bar{x}$. Letting $k \xrightarrow{\mathcal{N}} \infty$ in (5.17) and taking into account the convergence above and (5.19) bring us to the estimate

$$\|\nabla f(\bar{x})\|^2 \leq \beta \|\nabla f(\bar{x})\|^2.$$

This contradicts $\beta < \frac{1}{2}$ and $\|\nabla f(\bar{x})\| > 0$. Thus the origin is an accumulation point of $\{g^k\}$ as claimed.

Claim 2 Every accumulation point of $\{x^k\}$ is a stationary point of f .

Pick any accumulation point \bar{x} of $\{x^k\}$. Using Claim 1, the second inequality in (5.13), and Lemma 2.2 tells us that there exists an infinite set $J \subset \mathbb{N}$ such that

$$x^k \xrightarrow{J} \bar{x} \text{ and } g^k \xrightarrow{J} 0.$$

Take a neighborhood Ω of \bar{x} and $\Delta > 0$. Since \mathcal{G} is a local approximation of ∇f under condition (3.2), there exists $C > 0$ for which

$$\|\mathcal{G}(x, \delta) - \nabla f(x)\| \leq C\delta \text{ whenever } 0 < \delta \leq \Delta \text{ and } x \in \Omega. \quad (5.20)$$

Since $v_k \downarrow 0$ and $x^k \xrightarrow{J} \bar{x}$, we can select $K \in \mathbb{N}$ so that $v_k \leq \Delta$ and $x^k \in \Omega$ for all $k \geq K$, $k \in J$. This ensures together with (5.20) that

$$\|g^k - \nabla f(x^k)\| = \|\mathcal{G}(x^k, \min\{\delta_{k+1}, v_k\}) - \nabla f(x^k)\| \leq C v_k \text{ for all } k \geq K, k \in J.$$

Employing $g^k \xrightarrow{J} 0$ and $v_k \downarrow 0$ as above, we deduce that $\nabla f(x^k) \xrightarrow{J} 0$, and hence $\nabla f(\bar{x}) = 0$. Therefore, \bar{x} is a stationary point of f , which justifies (i).

Now we verify (ii) and (iii) simultaneously. It follows from (5.13) and $\tau_k \leq 1$ for all $k \in \mathbb{N}$ by the choice of τ_k in Step 3 of Algorithm 3 that

$$\sum_{k=1}^{\infty} \|x^{k+1} - x^k\|^2 = \sum_{k=1}^{\infty} \tau_k^2 \|g^k\|^2 \leq \bar{\tau} \sum_{k=1}^{\infty} \tau_k \|g^k\|^2 < \infty,$$

which implies that $\|x^{k+1} - x^k\| \rightarrow 0$. Then both assertions (ii) and (iii) follow from Lemma 2.3. \square

The next result establishes the global convergence with convergence rates of the iterates $\{x^k\}$ in Algorithm 3 under the KL property and the boundedness of $\{x^k\}$. We have already discussed the KL property in Remark 2.5. The boundedness of $\{x^k\}$ is also a standard assumption that appears in many works on gradient descent methods; see, e.g., [4, Theorem 4.1], [36, Theorem 1], and [46, Assumption 7].

Theorem 5.5 *Let $\{x^k\}$ be the sequence of iterates generated by Algorithm 3. Assuming that $\nabla f(x^k) \neq 0$ for all $k \in \mathbb{N}$ and that $\{x^k\}$ is bounded yields the assertions:*

(i) *If \bar{x} is an accumulation point of $\{x^k\}$ and f satisfies the KL property at \bar{x} , then $x^k \rightarrow \bar{x}$ as $k \rightarrow \infty$.*

(ii) *If in addition to (i), the KL property at \bar{x} is satisfied with $\psi(t) = Mt^q$ for some $M > 0$, $q \in [1/2, 1)$, then the following convergence rates are guaranteed:*

- *If $q = 1/2$, then the sequence $\{x^k\}$ converges linearly to \bar{x} .*
- *If $q \in (1/2, 1)$, then we have the estimate*

$$\|x^k - \bar{x}\| = \mathcal{O}(k^{-\frac{1-q}{2q-1}}).$$

Proof Let $\Omega := \{x^k\}$, and let $\Delta > 0$. Since \mathcal{G} is a local approximation of ∇f satisfying condition (3.2), there exists a positive number C such that

$$\|\mathcal{G}(x, \delta) - \nabla f(x)\| \leq C\delta \quad \text{whenever } x \in \Omega \text{ and } 0 < \delta \leq \Delta. \quad (5.21)$$

Select $K \in \mathbb{N}$ so that $v_k < \Delta$ for all $k \geq K$, which implies by (5.21) and the choice of g^k in Step 1 of Algorithm 3 the relationships

$$\begin{aligned} \|g^k - \nabla f(x^k)\| &= \|\mathcal{G}(x^k, \min\{\delta_{k+1}, v_k\}) - \nabla f(x^k)\| \\ &\leq C \min\{\delta_{k+1}, v_k\} \leq C\delta_{k+1} \quad \text{for all } k \geq K. \end{aligned} \quad (5.22)$$

We split the proof of the result into two parts by showing first that the sequences $\{C_k\}$ and $\{t_k^{\min}\}$ are constant after a finite number of iterations and verifying then the convergence of $\{x^k\}$ in (i) with the rates in (ii) by using Propositions 2.6 and 2.7.

Claim 3 *There exists $k_0 \in \mathbb{N}$ such that $C_k = C_{k_0}$ and $t_k^{\min} = t_{k_0}^{\min}$ for all $k \geq k_0$.*

Arguing by contradiction, suppose that such a number k_0 does not exist. By the construction of $\{C_k\}$ and $\{t_k^{\min}\}$ in Step 3 of Algorithm 3, we deduce that $C_k \uparrow \infty$ and $t_k^{\min} \downarrow 0$ as $k \rightarrow \infty$. Since Ω is bounded, Lemma 5.2 allows us to find $\bar{t} \in (0, 1)$ for which

$$\begin{aligned} f(x - tg) &\leq f(x) - \beta t \|g\|^2 \quad \text{whenever } x \in \Omega, \quad \|g - \nabla f(x)\| \leq \frac{1}{2} \|g\|, \\ &\text{and } t \in (0, \bar{t}]. \end{aligned} \quad (5.23)$$

Using the aforementioned properties of $\{C_k\}$ and $\{t_k^{\min}\}$, we get $N \geq K$ such that $C_k > C$ and $t_k^{\min} < \gamma \bar{t}$ for all $k \geq N$. Fix such a number k and then combine the

condition $\|g^k\| > \mu C_k \delta_{k+1}$ from (5.1) with $C_k > C$, $\mu > 2$, and (5.22). This gives us the inequalities

$$\|g^k\| > \mu C_k \delta_{k+1} \geq \mu C \delta_{k+1} \geq 2 \|g^k - \nabla f(x^k)\|,$$

which imply together with $x^k \in \Omega$ and (5.23) the estimate

$$f(x^k - t g^k) \leq f(x^k) - \beta t \|g^k\|^2 \text{ for all } t \in (0, \bar{t})$$

and thus tell us that $t_k > \gamma \bar{t} > t_k^{\min}$. Employing Step 3 of Algorithm 3 yields $t_{k+1}^{\min} = t_k^{\min}$. Since the latter holds whenever $k \geq N$, we conclude that the equality $t_k^{\min} = t_N^{\min}$ is satisfied for all $k \geq N$. This contradicts the condition $t_k^{\min} \downarrow 0$ as $k \rightarrow \infty$ and hence justifies the claimed assertion.

Claim 4 *All the assertions in (i) and (ii) are fulfilled.*

From Step 2 and Step 3 of Algorithm 3, we deduce that

$$f(x^k) - f(x^{k+1}) \geq \beta \tau_k \|g^k\|^2 \text{ for all } k \in \mathbb{N}. \quad (5.24)$$

Defining $N := \max\{K, k_0\}$ with k_0 taken from Claim 3 gives us the equalities

$$C_k = C_N \text{ and } t_k^{\min} = t_N^{\min} \text{ whenever } k \geq N. \quad (5.25)$$

Combining $C_k = C_N$ with (5.22) and $\|g^k\| > \mu C_k \delta_{k+1}$ from (5.1) ensures that

$$\begin{aligned} \|\nabla f(x^k)\| &\leq \|g^k\| + C \delta_{k+1} \\ &\leq \|g^k\| + \frac{C}{\mu C_N} \|g^k\| = \alpha \|g^k\| \text{ for all } k \geq N, \end{aligned} \quad (5.26)$$

where $\alpha := 1 + \frac{C}{\mu C_N}$. In addition, we have $t_{k+1}^{\min} = t_k^{\min} = t_N^{\min}$ in (5.25), which implies together with Step 3 of Algorithm 3 the relationships

$$\tau_k = t_k \geq t_k^{\min} = t_N^{\min} \text{ as } k \geq N \quad (5.27)$$

confirming the boundedness of $\{\tau_k\}$ from below. If the KL property of f holds at the accumulation point \bar{x} of $\{x^k\}$, it follows from Remark 2.8(i), (5.24), (5.26), and (5.27) that assumptions (H1) and (H2) in Proposition 2.6 hold. Thus $x^k \rightarrow \bar{x}$ as $k \rightarrow \infty$, which verifies (i).

Assume finally that the KL property at \bar{x} is satisfied with $\psi(t) = Mt^q$, $M > 0$, and $q \in [1/2, 1)$. The iterative procedure $x^{k+1} = x^k - \tau_k g^k$ in Step 4 of Algorithm 3 together with (5.27) and $g^k > 0$ from Step 1 therein tells us that $x^{k+1} \neq x^k$ for $k \geq N$. Combining this with (5.24), (5.26), and (5.27) verifies all the assumptions of Proposition 2.7 and therefore completes the proof of the theorem. \square

5.2 Dynamic step linesearch for noisy functions

In this subsection, we continue the study of problem (1.1) with the objective function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ of class $\mathcal{C}^{1,1}$. Similarly to Subsection 4.3, assume that only a *noisy approximation* $\phi(x) = f(x) + \xi(x)$ of f is available, where $\xi : \mathbb{R}^n \rightarrow \mathbb{R}$ is a noise function bounded by some known constant $\xi_f > 0$. Unlike Subsection 4.3, which considers $\mathcal{C}_L^{1,1}$ functions with an unknown noise level, we assume here that ξ_f is known, which may at first seem impractical or inefficient. However, if the noise is generated by independent and identically distributed random variables, we can employ an additional minor step to approximate the noise locally and use this approximation as a global noise level. This implementation will be discussed more carefully in the numerical experiments in Subsection 6.2. Considering only the forward finite difference given by

$$\tilde{\mathcal{G}}(x, \delta) = \frac{1}{\delta} \sum_{i=1}^n (\phi(x + \delta e_i) - \phi(x)) e_i \text{ for any } (x, \delta) \in \mathbb{R}^n \times (0, \infty), \quad (5.28)$$

we state the following noisy version of Proposition 3.4 that can be verified similarly.

Proposition 5.6 *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a \mathcal{C}^1 -smooth function such that ∇f is Lipschitz continuous on $\mathbb{B}(x, \Delta)$ with constant $\ell > 0$. Then the noisy forward finite difference (5.28) satisfies the error bound*

$$\|\tilde{\mathcal{G}}(x, \delta) - \nabla f(x)\| \leq \frac{\ell \sqrt{n} \delta}{2} + \frac{2\sqrt{n} \xi_f}{\delta} \text{ for all } \delta \in (0, \Delta]. \quad (5.29)$$

Now we design Derivative-Free Method with Dynamic Step Linesearch (DFD), which main feature is a *dynamic step linesearch* for determining both the *stepsize* and the *finite difference interval*.

Algorithm 4 (DFD for noisy functions).

Step 0 (initialization). Select an initial point $x^1 \in \mathbb{R}^n$, $\eta > 1$, and $L_1 > 0$. Set $k := 1$.

Step 1 (dynamic step linesearch). Find $i_k \in \mathbb{Z}$ with the smallest absolute value such that for $g^k := \tilde{\mathcal{G}}\left(x^k, \sqrt{\frac{4\xi_f}{\eta^{i_k} L_k}}\right)$ and $\tau_k = \frac{1}{\eta^{i_k} L_k}$, it holds that

$$\phi\left(x^k - \tau_k g^k\right) \leq \phi(x^k) - \frac{\tau_k}{9} \|g^k\|^2. \quad (5.30)$$

Step 2 (stepsize and parameters update). Set $x^{k+1} := x^k - \tau_k g^k$ and $L_{k+1} := \eta^{i_k} L_k$.

Remark 5.7 (Discussions on dynamic step linesearch)

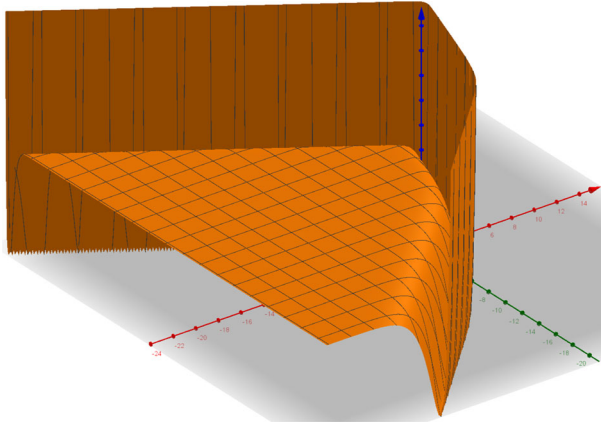


Fig. 1 Graph of $f(x, y) = (e^{2x+3y-1} + e^{3x-y} + e^{x-y-6} - 3)^2$

- In Step 1 of Algorithm 4, we employ dynamic step linesearch to find an approximation $L_{k+1} := \eta^{ik} L_k$ for the Lipschitz constant ℓ_k of ∇f locally around the current iterate x^k . Then we use L_{k+1} to determine both the stepsize $\tau_k = \frac{1}{L_{k+1}}$ and the finite difference interval $\sqrt{\frac{4\xi_f}{L_{k+1}}}$, which is the minimizer of the right-hand side of (5.29) with respect to δ when $\ell = L_{k+1}$. This idea also appeared in [20] for regularized Newton methods, where the finite difference interval is adaptively adjusted and used for determining the cubic regularization parameter.
- The attempt of using dynamic step linesearch in order to adjust the stepsize is not new as it has been already employed in [24, 60] for gradient descent methods with the exact gradient and in [7, 8] for derivative-free optimization methods. However, employing this procedure to determine both *stepsize* and *finite difference interval* is suggested, to the best of our knowledge, in the present paper for the first time.
- The dynamic step linesearch procedure plays an important role not only in our convergence analysis but also in practical modeling. *Theoretically*, condition (5.30) is necessary for the value $\eta^{ik} L_k$ in Step 1 being a good approximation of the Lipschitz constant ℓ_k of ∇f locally around the reference iterate x^k , which is confirmed by Proposition 5.9. *Numerically*, by automatically approximating the local Lipschitz constant of the gradient, our DFD has a better performance in comparison with other finite-difference-based algorithms for noisy $C^{1,1}$ functions with complex structures. To illustrate this claim, we consider the function $f(x, y) := (e^{2x+3y-1} + e^{3x-y} + e^{x-y-6} - 3)^2$ of two variables with the graph in \mathbb{R}^3 depicted below.

This function is inspired by a univariate function in [63, Section 4.1.2]. It provides a challenging example for finite-difference-based methods in both cases of approximate gradients and finding minimizers. This is because f has very small first- and second-order derivatives at points belonging to most of the second and third parts of the plane. In addition, minimizing the function is more challenging in the context of derivative-free optimization because the approximate gradient

Table 1 Stepsize and finite difference interval selections

Method	IMFIL	RG	GDD (Ada)	L-BFGS (Ada)	DF-backtracking	DFD
Stepsize	Backtracking	GS	Dynamic	Armijo + Wolfe	Backtracking	Dynamic
FD interval	Decreasing	GS	Adaptive	Adaptive	Backtracking	Dynamic

obtained from finite differences may be unreliable, as the function values change rapidly between flat and sharp regions. As in the context of noisy DFO, we manually inject into f uniformly distributed stochastic noises with different levels. The plots below show the trajectories of iterates generated by our DFD (Algorithm 4) and the following algorithms:

- IMFIL: The implicit filtering algorithm [26, Algorithm 2.2] with the forward finite difference.
- RG: The random gradient-free method [52, Section 5].
- L-BGFS (Ada): The noise-tolerant quasi-Newton algorithm [64, Algorithm 2.1], where the gradient is approximated by the forward finite difference with the adaptive finite difference interval estimation from [63, Algorithm 2.1].
- GDD (Ada): Gradient descent with dynamic step linesearch, where the gradient is approximated by the forward finite difference with the adaptive finite difference interval estimation from [63, Algorithm 2.1].
- DF-backtracking: A modified version of our basic DFD, where the dynamic step linesearch is replaced by the standard backtracking linesearch, i.e., the condition $i \in \mathbb{Z}$ is replaced by $i \in \mathbb{N}$ in Step 1 of Algorithm 4.

In the algorithms above, only our DFD method (Algorithm 4) uses the dynamic step linesearch to determine both *stepsize* and *finite difference interval*. The selections of stepsize and finite difference interval for each method are listed in Table 1, where GS in the selections of RG means grid search. Details for the settings of the algorithms and additional numerical results on this experiment can be found in Appendix A, where different noise levels are addressed.

It can be observed from Figure 2 addressing the noise level 0.01 that only the last points (red stars) generated by our DFD method successfully identify the minimum region (depicted in dark blue) regardless of the choice of initial points (blue circles). L-BFGS (Ada) locates the minimum region in the only case when the initial point is $(-4, -4)$. Other algorithms including RG, IMFIL, GDD (Ada), and DF-backtracking perform even worse since they remain stuck at the initial points in all the scenarios. Additional graphs in Appendix A also show that these results are stable with respect to different levels of noise ranging from 1 or 10^{-3} . The failure of GDD (Ada) and DF-backtracking emphasizes the crucial role of using dynamic step linesearch to determine both the stepsize and the finite difference interval in the construction of DFD.

The rest of this subsection is devoted to deriving the fundamental convergence properties of Algorithm 4 for noisy smooth functions. We begin with a simple albeit useful lemma about the *optimal local Lipschitz constant* of the gradient of a $\mathcal{C}^{1,1}$ function.

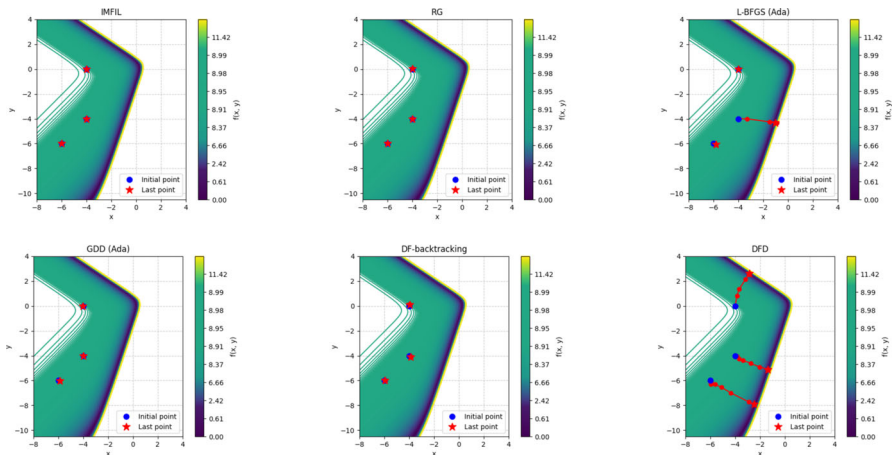


Fig. 2 Finite-difference-based methods on minimizing a $C^{1,1}$ function with complex structure

Lemma 5.8 *Let $\Omega \subset \mathbb{R}^n$ be a nonempty bounded set. Then for any $\xi > 0$, there exists some $\ell > 0$ such that ℓ is the Lipschitz constant of ∇f on $\bigcup_{x \in \Omega} \mathbb{B}(x, \delta_x)$, where*

$$\delta_x := \max \left\{ \frac{3}{2\ell} \|\nabla f(x)\|, \sqrt{\frac{4\xi}{\ell}} \right\}.$$

Proof Define the number $M := \sup \{ \|\nabla f(x)\| \mid x \in \Omega \} \in \mathbb{R}$ and deduce from the assumed $C^{1,1}$ property of f that f is Lipschitz continuous on the set $\bigcup_{x \in \Omega} \mathbb{B}(x, \max \{ \frac{3}{2}M, 2\sqrt{\xi} \})$ with some Lipschitz constant $L > 0$. Denoting $\ell := \max \{1, L\}$, we get

$$\delta_x = \max \left\{ \frac{3}{2\ell} \|\nabla f(x)\|, \sqrt{\frac{4\xi}{\ell}} \right\} \leq \max \left\{ \frac{3}{2}M, 2\sqrt{\xi} \right\} \quad \text{for all } x \in \Omega.$$

This tells us that ∇f is Lipschitz continuous with the constant ℓ on $\bigcup_{x \in \Omega} \mathbb{B}(x, \delta_x)$ as claimed. \square

The next result plays a crucial technical role in deriving the convergence properties in what follows.

Proposition 5.9 *Let $\ell > 0$ and $x \in \mathbb{R}^n$ be such that ∇f is Lipschitz continuous with some constant $\ell > 0$ on $\mathbb{B}(x, \max \{ \frac{3}{2\ell} \|\nabla f(x)\|, \sqrt{\frac{4\xi_f}{\ell}} \})$, and let $\tilde{\ell} > 0$, $i \in \mathbb{Z}$, $\eta > 1$ be selected so that $\ell \in (\eta^{i-1}\tilde{\ell}, \eta^i\tilde{\ell}]$. Define $g \in \mathbb{R}^n$ and $\tau > 0$ by*

$$g := \tilde{\mathcal{G}} \left(x, \sqrt{\frac{4\xi_f}{\eta^i\tilde{\ell}}} \right) \quad \text{and} \quad \tau := \frac{1}{\eta^i\tilde{\ell}},$$

where $\tilde{\mathcal{G}}$ is taken from (5.28). If $\|\nabla f(x)\| \geq 8\sqrt{\ell\eta n\xi_f}$, then we have the estimates

- (i) $f(x - \tau g) \leq f(x) - \frac{3\tau}{32} \|\nabla f(x)\|^2$,
(ii) $\phi(x - \tau g) \leq \phi(x) - \frac{\tau}{9} \|g\|^2$.

Proof Since ℓ is the Lipschitz constant of ∇f on $\mathbb{B}\left(x, \sqrt{\frac{4\xi_f}{\ell}}\right)$, we deduce from Proposition 5.28 that

$$\|\tilde{\mathcal{G}}(x, \delta) - \nabla f(x)\| \leq \frac{\ell\sqrt{n}\delta}{2} + \frac{2\sqrt{n}\xi_f}{\delta} \text{ for all } \delta \in \left(0, \sqrt{\frac{4\xi_f}{\ell}}\right]. \quad (5.31)$$

Combining this with $g = \tilde{\mathcal{G}}\left(x, \sqrt{\frac{4\xi_f}{\eta^i L}}\right)$ and $\ell \leq \eta^i L$ tells us that

$$\begin{aligned} \|g - \nabla f(x)\| &\leq \frac{\ell\sqrt{n}}{2} \sqrt{\frac{4\xi_f}{\eta^i L}} + 2\sqrt{n}\xi_f \sqrt{\frac{\eta^i \ell}{4\xi_f}} \\ &\leq \frac{\eta^i \tilde{\ell}\sqrt{n}}{2} \sqrt{\frac{4\xi_f}{\eta^i L}} + 2\sqrt{n}\xi_f \sqrt{\frac{\eta^i L}{4\xi_f}} = 2\sqrt{\eta^i \tilde{\ell} n \xi_f}. \end{aligned}$$

Using the triangle inequality and $\eta^{i-1}\tilde{\ell} < \ell$ yields

$$\begin{aligned} \|g\| &\geq \|\nabla f(x)\| - \|g - \nabla f(x)\| \\ &\geq 8\sqrt{\eta \ell n \xi_f} - 2\sqrt{\eta^i \tilde{\ell} n \xi_f} \\ &> 6\sqrt{\eta^i \tilde{\ell} n \xi_f} \geq 3\|g - \nabla f(x)\|, \end{aligned} \quad (5.32)$$

which being combined with the Cauchy-Schwarz inequality ensures that

$$\begin{aligned} \langle \nabla f(x), g \rangle &= \langle \nabla f(x) - g, g \rangle + \|g\|^2 \\ &\geq -\|\nabla f(x) - g\| \|g\| + \|g\|^2 \geq \frac{2}{3} \|g\|^2. \end{aligned}$$

Thus we arrive at $\|g\| \leq \frac{3}{2} \|\nabla f(x)\|$ implying together with $\tau = \frac{1}{\eta^i L} \leq \frac{1}{\ell}$ that

$$x - \tau g \in \mathbb{B}\left(x, \frac{3}{2\eta^i L} \|\nabla f(x)\|\right) \subset \mathbb{B}\left(x, \frac{3}{2\ell} \|\nabla f(x)\|\right).$$

By the Lipschitz continuity of ∇f with constant ℓ on the ball above and Lemma 2.1, we get

$$\begin{aligned} f(x - \tau g) &\leq f(x) + \langle x - \tau g - x, \nabla f(x) \rangle + \frac{\ell}{2} \|x - \tau g - x\|^2 \\ &= f(x) - \tau \langle g, \nabla f(x) \rangle + \frac{\ell\tau^2}{2} \|g\|^2 \\ &\leq f(x) - \frac{2\tau}{3} \|g\|^2 + \frac{\tau}{2} \|g\|^2 = f(x) - \frac{\tau}{6} \|g\|^2. \end{aligned} \quad (5.33)$$

It also follows from (5.32) that

$$\|g\| \geq \|\nabla f(x)\| - \|g - \nabla f(x)\| \geq \|\nabla f(x)\| - \frac{1}{3} \|g\|,$$

which yields $\|\nabla f(x)\| \leq \frac{4}{3} \|g\|$ and, being combined with (5.33), verifies (i).

(ii) Using (5.33) and the construction of the noisy approximation ϕ gives us the estimate

$$\phi(x - \tau g) \leq \phi(x) - \frac{\tau}{6} \|g\|^2 + 2\xi_f, \quad (5.34)$$

which implies together with $\eta^{i-1}\tilde{\ell} < \ell$, $n \geq 1$, and $\|\nabla f(x)\| \geq 8\sqrt{\ell\eta n\xi_f}$ that

$$\frac{\tau}{18} \|g\|^2 \geq \frac{1}{18\eta^i\tilde{\ell}} \frac{9}{16} \|\nabla f(x)\|^2 \geq \frac{2}{64\eta^i\tilde{\ell}} 64\eta\ell n\xi_f \geq 2\xi_f.$$

Combining the latter with (5.34) leads us to the conclusion in (ii) and thus completes the proof. \square

Similarly to Subsection 4.3, we say that Step 1 of Algorithm 4 is *successful* if the integer number i_k is found, and *unsuccessful* otherwise. It follows directly from Proposition 5.9 that Step 1 of Algorithm 4 is successful whenever $\|\nabla f(x^k)\|$ is not near 0 as stated below.

Corollary 5.10 *At the k^{th} iteration of Algorithm 4, let ℓ_k be such that ∇f is Lipschitz continuous on $\mathbb{B}(x^k, \max\{\frac{3}{2\ell_k} \|\nabla f(x^k)\|, \sqrt{\frac{4\xi_f}{\ell_k}}\})$ with some constant $\ell_k > 0$. If the condition*

$$\|\nabla f(x^k)\| \geq 8\sqrt{\ell_k\eta n\xi_f} \quad (5.35)$$

is satisfied, then Step 1 of Algorithm 4 is successful.

Remark 5.11 The result above is one of the crucial findings that illustrate behavior of Algorithm 4. It shows that as long as $\nabla f(x^k)$ is not small relative to the *noise* and the *local Lipschitz constant* of ∇f around x^k , the algorithm always makes significant progress. This also explains the success of DFD in the experiment presented in Figure 2. In the flat regions where the gradient magnitude is small, the local Lipschitz constant of the gradient is small as well, which ensures that condition (5.35) remains valid. Consequently, the iterative sequence attempts to move out of such regions.

Employing the obtained corollary, we arrive at the next proposition, which is useful in the proof of the main convergence results below.

Proposition 5.12 *At some k^{th} iteration of Algorithm 4, let $L > 0$ be such that ∇f is Lipschitz continuous with constant L on $\mathbb{B}(x^k, \max\{\frac{3}{2L} \|\nabla f(x^k)\|, \sqrt{\frac{4\xi_f}{L}}\})$ and assume that*

$$\|\nabla f(x^k)\| \geq 8\sqrt{L\eta n\xi_f}.$$

The following assertions hold:

- (i) If $L_k < \eta L$ then $L_{k+1} < \eta L$.
- (ii) If $L_k \geq L$ then $L_{k+1} \geq L$.
- (iii) If $L_k \in [L, \eta L)$ then $L_{k+1} = L_k$.

Proof (i) By the construction of $\{L_k\}$, we find $m \in \mathbb{Z}$ such that $L_{k+1} = \eta^m L_k$. If $m \leq 0$, then $L_{k+1} \leq L_k < \eta L$, and so we assume that $m > 0$. Then the exit condition in Step 1 of Algorithm 4 yields

$$\phi\left(x^k - \frac{1}{\eta^i L_k} g_i^k\right) > \phi(x^k) - \frac{1}{\eta^i L_k} \|g_i^k\|^2 \text{ for all } i \in \{0, \dots, m-1\}, \quad (5.36)$$

where $g_i^k := \tilde{G}(x^k, \sqrt{\frac{4\varepsilon_f}{\eta^i L_k}})$. Observe that condition (5.35) holds for $\ell_k = L$, which is a Lipschitz constant of ∇f on $\mathbb{B}(x^k, \max\{\frac{3}{2L} \|\nabla f(x^k)\|, \sqrt{\frac{4\varepsilon_f}{L}}\})$ by the assumptions made. Combining this with Corollary 5.10 and estimate (5.36), we deduce that $\eta^i L_k \notin [L, \eta L)$ for all $i \in \{0, \dots, m-1\}$. This fact together with $L_k < \eta L$ tells us that $L_{k+1} = \eta^m L_k < \eta L$.

(ii) By the construction of $\{L_k\}$, we find some $m \in \mathbb{Z}$ such that $L_{k+1} = \eta^m L_k$. If $m \geq 0$, then $L_{k+1} \geq L_k \geq L$, and so we assume that $m < 0$. Then the exit condition in Step 1 of Algorithm 4 yields

$$\phi\left(x^k - \frac{1}{\eta^i L_k} g_i^k\right) > \phi(x^k) - \frac{1}{\eta^i L_k} \|g_i^k\|^2 \text{ for all } i \in \{0, -1, \dots, m+1\}, \quad (5.37)$$

where $g_i^k := \tilde{G}(x^k, \sqrt{\frac{4\varepsilon_f}{\eta^i L_k}})$. Observe that condition (5.35) holds for $\ell_k = L$, which is also a Lipschitz constant of ∇f on $\mathbb{B}(x^k, \max\{\frac{3}{2L} \|\nabla f(x^k)\|, \sqrt{\frac{4\varepsilon_f}{L}}\})$ by the assumptions made. Combining this with Corollary 5.10 and (5.37), we get that $\eta^i L_k \notin [L, \eta L)$ for all $i \in \{0, -1, \dots, m+1\}$. This fact together with $L_k \geq L$ verifies that $L_{k+1} = \eta^m L_k \geq L$. \square

Now we are in a position to derive convergence properties of DFD from Algorithm 4. Consider first the case where at some K^{th} iteration, Step 1 of Algorithm 4 is not successful, i.e., we cannot find $i_K \in \mathbb{Z}$ that ensures the descent condition (5.30). Then Corollary 5.10 tells us that $\|\nabla f(x^K)\| < 8\sqrt{\ell_K \eta n \xi_f}$, where ℓ_K is a Lipschitz constant of ∇f around x^K . In this case, Algorithm 4 finds a point near a stationary one after a finite number of iteration. In practice, to avoid the process of finding i_k in Step 1 of Algorithm 4 from running infinitely to cause a computational error, the users can add a lower bound sufficiently small and an upper bound sufficiently large for i_k in the loop.

The main theorem of this section concerns the case where Step 1 of Algorithm 4 is successful for all $k \in \mathbb{N}$. In this scenario, we can find a point near a stationary one, along the sequence of iterates generated by the algorithm, if just one of the Lipschitz approximations is appropriate.

Theorem 5.13 Assume that Step 1 of Algorithm 4 is successful for all $k \in \mathbb{N}$ and that there exists $L > 0$ such that ∇f is Lipschitz continuous with constant L on $\bigcup_{k=1}^{\infty} \mathbb{B}(x^k, \max\{\frac{3}{2L} \|\nabla f(x^k)\|, \sqrt{\frac{4\xi_f}{L}}\})$. If $\inf_{k \in \mathbb{N}} f(x^k) > -\infty$ and for some $K \in \mathbb{N}$ we have $L_K \in [L, \eta L)$, then the following assertions hold:

(i) There exists $N \in \mathbb{N}$ for which

$$\|\nabla f(x^N)\| < 8\sqrt{L\eta\xi_f}. \quad (5.38)$$

(ii) Assume in addition that f has a global minimizer with the minimum value f^* , that $f(x^K) > f^*$, and that f satisfies the Polyak-Łojasiewicz inequality with some constant $\mu > 0$, i.e.,

$$\mu(f(x) - f^*) \leq \frac{1}{2} \|\nabla f(x)\|^2 \quad \text{for all } x \in \mathbb{R}^n. \quad (5.39)$$

Then the number N from (5.38) admits the upper estimate

$$N \leq \max \left\{ 1 + K, 1 + K + \log_{1 - \frac{3\mu}{16\eta L}} \left(\frac{32\eta\xi_f}{f(x^K) - f^*} \right) \right\}. \quad (5.40)$$

Proof (i) Assume on the contrary that $\|\nabla f(x^k)\| \geq 8\sqrt{L\eta\xi_f}$ as $k \in \mathbb{N}$. It follows from Proposition 5.12 and $L_K \in [L, \eta L)$ that $L_{k+1} = L_k$ whenever $k \geq K$. Using Proposition 5.9(i) with

$$\ell := L, \quad x := x^k, \quad \tilde{\ell} := L_k, \quad \text{and } i := 0,$$

we get the relationship below between two subsequent iterations

$$f(x^{k+1}) \leq f(x^k) - \frac{3}{32L_K} \|\nabla f(x^k)\| \quad \text{whenever } k \geq K,$$

which tells us that $\{f(x^k)\}$ is a strictly decreasing sequence. By $\inf_{k \in \mathbb{N}} f(x^k) > -\infty$, this sequence is convergent, and hence $\|\nabla f(x^k)\| \rightarrow 0$ as $k \rightarrow \infty$. We arrive at a contradiction with $\|\nabla f(x^k)\| \geq 8\sqrt{L\eta\xi_f}$ for all $k \in \mathbb{N}$, and thus justify (5.38) in (i).

To verify now assertion (ii), let N be the first iteration for which (5.38) holds, i.e.,

$$\|\nabla f(x^k)\| \geq 8\sqrt{L\eta\xi_f} \quad \text{for } k \in \{1, \dots, N-1\}.$$

If $N \leq K+1$, estimate (5.40) is obviously satisfied, and thus we suppose that $N > K+1$. It follows from Proposition 5.12 that $L_k = L_K \in [L, \eta L)$ for all $k \in \{K, \dots, N-1\}$. Fixing such a number k and employing Proposition 5.9(i) for

$$\ell := L, \quad x := x^k, \quad \tilde{\ell} := L_k, \quad \text{and } i := 0$$

clearly bring us to the estimates

$$f(x^{k+1}) \leq f(x^k) - \frac{3}{32L_K} \left\| \nabla f(x^k) \right\|^2 \leq f(x^k) - \frac{3}{32\eta L} \left\| \nabla f(x^k) \right\|^2.$$

Combining this with the Polyak-Łojasiewicz inequality from (5.39), we obtain the condition

$$f(x^{k+1}) \leq f(x^k) - \frac{3\mu}{16\eta L} (f(x^k) - f^*),$$

which can be equivalently rewritten as

$$f(x^{k+1}) - f^* \leq \left(1 - \frac{3\mu}{16\eta L}\right) (f(x^k) - f^*).$$

Using the latter condition for $k = K, K + 1, \dots, N - 2$ gives us

$$f(x^{N-1}) - f^* \leq \left(1 - \frac{3\mu}{16\eta L}\right)^{N-1-K} (f(x^K) - f^*). \quad (5.41)$$

Since ∇f is Lipschitz continuous on $\overline{\mathbb{B}}(x^{N-1}, \frac{1}{L} \nabla f(x^{N-1}))$ with constant L , Lemma 2.1 yields

$$\begin{aligned} f^* &\leq f\left(x^{N-1} - \frac{1}{L} \nabla f(x^{N-1})\right) \\ &\leq f(x^{N-1}) + \left\langle x^{N-1} - \frac{1}{L} \nabla f(x^{N-1}) - x^{N-1}, \nabla f(x^{N-1}) \right\rangle \\ &\quad + \frac{L}{2} \left\| x^{N-1} - \frac{1}{L} \nabla f(x^{N-1}) - x^{N-1} \right\|^2 \\ &= f(x^{N-1}) - \frac{1}{L} \left\| \nabla f(x^{N-1}) \right\|^2 + \frac{1}{2L} \left\| \nabla f(x^{N-1}) \right\|^2, \end{aligned}$$

which ensures in turn the fulfillment of

$$f(x^{N-1}) - f^* \geq \frac{1}{2L} \left\| \nabla f(x^{N-1}) \right\|^2 \geq \frac{1}{2L} 64L\eta n\xi_f = 32\eta n\xi_f.$$

Combining the obtained estimates with (5.41) tells us that

$$32\eta n\xi_f \leq \left(1 - \frac{3\mu}{16\eta L}\right)^{N-1-K} (f(x^K) - f^*),$$

and thus verifies the claimed conclusion (5.40). \square

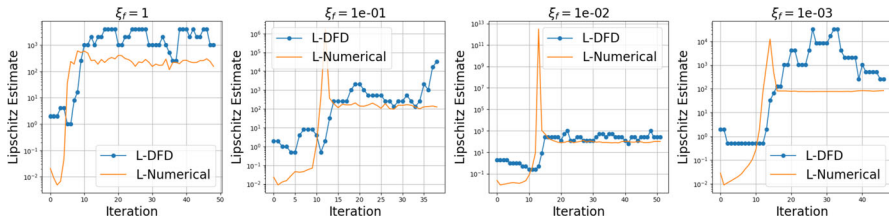


Fig. 3 $\{L_k\}$ generated by DFD approximates local Lipschitz constant of ∇f around iterates

Remark 5.14

- The existence of the constant L in the assumptions of Theorem 5.13 is guaranteed under the fulfillment of either one of the following conditions:
 - The objective function f is of class $\mathcal{C}_L^{1,1}$.
 - The level set $\{x \mid f(x) \leq f(x^1) + 2\xi_f\}$ is bounded. Indeed, it follows from (5.30) that the sets

$$\{x^k\} \subset \{x \in \mathbb{R}^n \mid \phi(x) \leq \phi(x^1)\} \subset \{x \in \mathbb{R}^n \mid f(x) \leq f(x^1) + 2\xi_f\}$$

are bounded as well. Combining the latter observation with Proposition 5.8 for $\Omega := \{x^k\}$ verifies the existence of the Lipschitz constant L .

- The assumption on the existence of $L_K \in [L, \eta L]$ in Theorem 5.13 is motivated by the construction Step 1 of Algorithm 4 and by our analysis conducted in Proposition 5.9, which shows that L_k being close to the true local Lipschitz constant of ∇f around x^k is related to the success of Step 1. In fact, this phenomenon also appears in practice as can be seen in Figure 3. The figure illustrates the sequence $\{L_k\}$ constructed in Algorithm 4 in comparison with the numerical Lipschitz constant of ∇f around x^k evaluated by

$$\max_{1 \leq i < j \leq 100} \frac{\|\nabla f(x_i^k) - \nabla f(x_j^k)\|}{\|x_i^k - x_j^k\|}, \text{ where } x_i^k \in \mathbb{B}\left(x^k, \frac{1}{10} \|\nabla f(x^k)\|\right)$$

is chosen uniformly.

The figure also demonstrates that Algorithm 4 automatically approximates the local Lipschitz constant of ∇f , which is in fact key to driving good numerical performance.

However, we believe that rigorously ensuring the assumption $L_k \in [L, \eta L]$ requires significantly more efforts including further investigation into the geometry of specific problems and the imposition of additional assumptions on the noise. Therefore, we would defer this analysis to a future research, where this algorithm is considered solely in more details.

6 Numerical experiments

In this section, we present numerical experiments demonstrating the efficiency of our methods in solving derivative-free optimization problems with and without the presence of noise. This section is split into two subsections addressing different noise levels: small noise, which also includes the noiseless case, and large noise. For each type of the noise level, we compare the performance of our newly developed methods with various well-known algorithms to ensure the diversity of the numerical experiments. In total, 786 test problems and 10 algorithms are considered in what follows.

6.1 Finite-difference-based algorithms for functions with small noise

Here we compare the performance of our DFC (Algorithm 1) and DFB (Algorithm 3) methods with other finite-difference-based algorithms to minimize smooth (convex and nonconvex) functions either without noise, or with small noise. The results in this subsection suggest that, in addition to the theoretical guarantees, our methods are more robust than the standard implementations of gradient descent methods with a *constant/backtracking stepsize* and with finite difference gradient for a *fixed finite difference interval*. The presented results also confirm the practicality of DFC and DFB methods in comparison with other well-known algorithms as in [26, 52].

6.1.1 Experiments with $\mathcal{C}_L^{1,1}$ functions

The first part of the subsection compares the performance of our DFC method using forward finite differences with some other well-known derivative-free methods for minimizing $\mathcal{C}_L^{1,1}$ functions. Since our DFC method is of the gradient descent type, we choose the set of testing algorithms as follows:

- (i) GDC (fixed), i.e., the standard gradient descent with a *constant stepsize* and gradients obtained from forward finite differences with a *fixed finite difference interval*.
- (ii) GD-ada, a variant of DFC with the stepsize being update by the rule in [6, Algorithm 2.2].
- (iii) IMFIL, i.e., the implicit filtering algorithm with forward finite differences [26].
- (iv) RG, i.e., a random gradient-free algorithm for smooth optimization proposed in [52].

The testing objective functions f are chosen as follows.

1. *Least-square (LS) regression*: $f(x) := \|Ax - b\|^2$, where A is an $n \times n$ matrix and $b \in \mathbb{R}^n$.
2. *A smooth nonconvex (NC) objective*: $f(x) := \sum_{i=1}^n \log(1 + (Ax - b)_i^2)$, where A is an $n \times n$ matrix and b is a vector in \mathbb{R}^n . This problem is considered in [59, Section 5.5] and [42, Section 4] with a nonsmooth term added to the objective function.

Random datasets are generated with different sizes for the testing purpose. To be more specific, an $n \times n$ matrix A and a vector $b \in \mathbb{R}^n$ are generated randomly with i.i.d. (independent and identically distributed) standard Gaussian entries. The dimension n is chosen from the set $\{10i, i = 1, \dots, 20\}$. We consider two types of noise in this

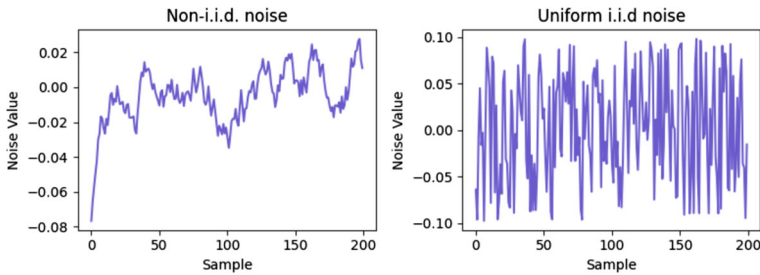


Fig. 4 Illustration for different types of noise

numerical experiment. The first type injects a uniformly distributed random noise with level $\xi_f \geq 0$, i.e., $\xi(x) \sim U(-\xi_f, \xi_f)$, into the function f and assumes the access only to $\phi(x) := f(x) + \xi(x)$ for all objective functions. In the second approach to create non-iid noise, let ε be a random array of length $200n$ constructed as follows:

$$\varepsilon_1 \sim U(-\xi_f, \xi_f), \quad \varepsilon_{k+1} := 0.9\varepsilon_k + 0.1\xi_k,$$

where $\xi_k \sim U(-\xi_f, \xi_f)$ for all $k \in \mathbb{N}$. Then for each evaluation at $x \in \mathbb{R}^n$, we set $\phi(x) := f(x) + \varepsilon_i$, where i is drawn randomly from $\{1, 2, \dots, 200n\}$. An illustration for these two types of noise with the same noise level $\xi_f = 0.1$ is presented below

The noise level is chosen from the set $\xi_f \in \{10^i, i = -9, \dots, -4\}$. The initial points are chosen as the zero vector for all the tests and algorithms. We also assume that the noise level is unknown in these numerical experiments. For that reason, the settings for DFC and GDC (fixed) are chosen as follows:

- DFC, GD-ada: The initial finite difference interval is $\delta_1 = 10^{-2}$. Other parameters are chosen as: $L_1 = n, \mu = 2.5, r = 2, \kappa = \sqrt{n}/2, \theta = 0.5$. The illustration in Figure 5 shows that L_1 and δ_1 are just loose lower bounds and loose upper bounds of the true Lipschitz constant and the optimal finite difference interval, respectively. This further emphasizes the practicality of the methods.
- GDC (fixed): The finite difference interval is chosen as $\delta = 10^{-8}$ for the noiseless case and $\delta = 2\sqrt{\xi_f}$ for the noisy case, which is of the same order as the optimal finite difference interval. Note that the latter selection is for testing purposes only since GDC (fixed) does not perform well with $\delta = 10^{-8}$ in the presence of noise. Of course, when the noise level is unknown, choosing a good finite difference interval for GDC (fixed) is not an easy task. To ensure a fair comparison, the stepsize of GDC (fixed) is chosen by a grid search on the set $\{\frac{1}{n}, \frac{0.2}{n}, \frac{0.1}{n}\}$, where n is the dimension of the problem.

The setting of IMFIL is similar to the one given in Appendix A. The setting of RG is also similar to that in Appendix A, except that the approximate Lipschitz constant is chosen by a grid search on $\{n, 5n, 10n\}$, where n is the dimension of the problem, to ensure a fair comparison. All the methods are executed until they reach the maximum number of function evaluations of $200n$.

In order to illustrate the performance of the algorithms, we use the performance profiles [19] with the measure $f_p^s - f_p^*$, where f_p^s is the function value obtained by

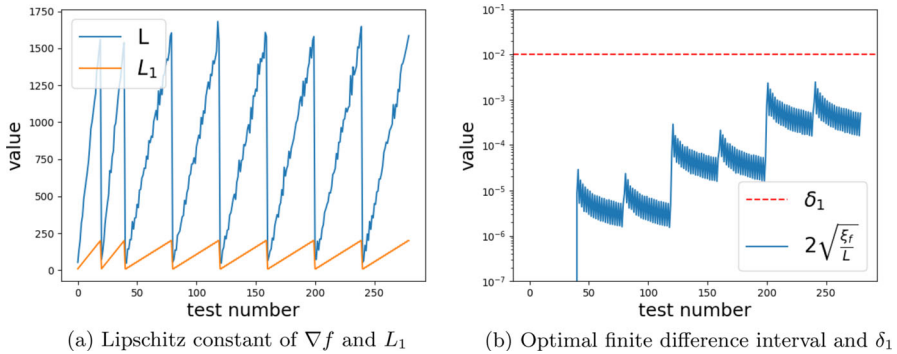


Fig. 5 Data for objective functions and initializations of DFC, GD-ada

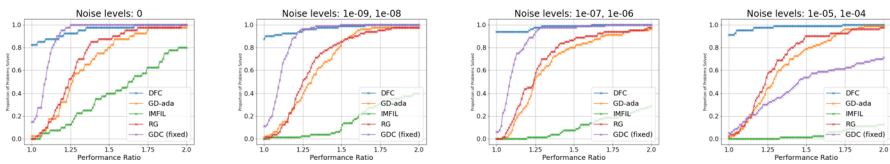


Fig. 6 Performance profiles of methods solving $C_L^{1,1}$ problems with uniform i.i.d. noise

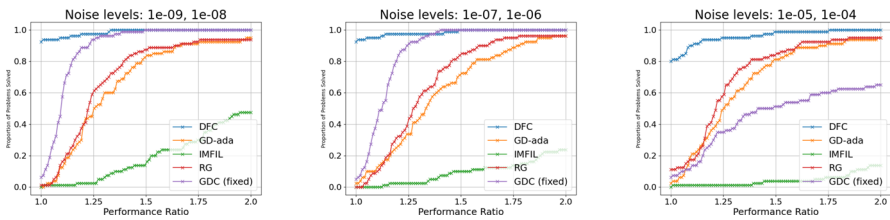


Fig. 7 Performance profiles of methods solving $C_L^{1,1}$ problems with non-i.i.d. noise

method s for problem p , and where f_p^* is the optimal value of the problem p . To be more specific, we assume that the set of problem tests is P . For each method s , we plot the graph of the function

$$\rho_s(\tau) := \frac{1}{|P|} \left| \left\{ p \in P \mid \frac{f_p^s - f_p^*}{f_p^{\text{best}} - f_p^*} \leq \tau \right\} \right| \text{ for } \tau \geq 1,$$

where $|P|$ is the size of P , and where f_p^{best} is the smallest function value obtained by all the methods in problem p . For example, $\rho_s(1)$ represents the percentage of problems where the method s performs the best. Due to the structure of the problems, f_p^* is always chosen to be 0. The results for different noise types and levels are presented in Figures 6 and 7. It can be seen that DFC performs the best in most tests. The robustness of DFC is also good for most selections of performance ratios and is increasing when the noise level is increasing.

Table 2 A set of unconstrained problems from CUTEst

Problem	n	Problem	n	Problem	n	Problem	n
ALLINITU	4	DIXMAANB	90	HIMMELBG	2	SPARSINE	100
ARWHEAD	100	DQRTIC	10	HIMMELBH	2	TOINTGSS	50
BARD	3	ENGVAL1	50	HUMPS	2	TOINTGSS	100
BDQRTIC	100	ENGVAL1	100	LOGHAIRY	2	TQUARTIC	100
BOX3	3	FLETBV3M	10	NCB20B	100	TRIDIA	100
BOXPOWER	100	FLETBV3M	100	NONDIA	100	VARDIM	10
BRKMCC	2	FLETGBV2	10	NONDQUAR	100	VAREIGVL	50
BROWNAL	100	FLETGBV3	10	PENALTY3	50	VAREIGVL	100
COSINE	10	FLETGBV3	100	POWELLSG	4	WOODS	100
CRAGGLVY	4	FLETCHCR	100	ROSENBRTU	2	ZANGWIL2	2
CURLY30	100	GULF	3	SENSORS	3		
DIXMAANB	15	HIMMELBCLS	2	SISSER	2		

6.1.2 Experiments with $\mathcal{C}^{1,1}$ Functions

In this subsection, we illustrate the performance of DFB method, i.e., Algorithm 3 with forward finite differences on a subset of CUTEst problems [23, 27] with the details given in Table 2. We also inject uniformly distributed stochastic noise as before, with the noise level ξ_f is either 0, or is chosen from the set $\{10^i, i = -9, \dots, -4\}$ while being unknown to the tested algorithms. In addition to DFB, the methods considered in this numerical experiment are IMFIL, RG with the same setting as in Subsection 6.1.1, and GDB (fixed), i.e., the standard gradient descent method with *backtracking stepsize*, where the approximate gradient is obtained from the forward finite difference with a *fixed finite difference interval*.

The settings for DFB and GDB (fixed) are chosen as follows:

- DFB: The initial finite difference interval $\delta_1 = 10^{-2}$. Other parameters are chosen as: $\theta = 0.5$, $\mu = 2.1$, $\eta = 2$, $\beta = 0.1$, $\gamma = 0.5$, $C_1 = \frac{\sqrt{n}}{2}$, $t_1^{\min} = 10^{-6}$, $\bar{\tau} = 1$, $v_k = 1/k$.
- GDB (fixed): The finite difference interval is chosen as $\delta = 10^{-8}$ for the noiseless case and $\delta = 2\sqrt{\xi_f}$ for the noisy case, similarly to the selection in GDC (fixed) in previous numerical experiments. The linesearch reduction factor is 0.5, the linesearch constant is 0.1, and the lower bound of the linesearch stepsize is 10^{-10} .

All the methods are executed until they reach the maximum number of function evaluations of $200 \cdot n$. Similarly to the previous experiments, the results here are illustrated by the performance profiles with the same measure as in Subsection 6.1.1. Since the exact optimal value is unknown, we approximate it by running DFB and Powell algorithm from SciPy library [66] with the maximum number of function evaluations of $400 \cdot n$ on the noiseless function. The performance profiles with different levels are presented in Figure 8 showing that DFB achieves the best performance for most of the performance ratios.

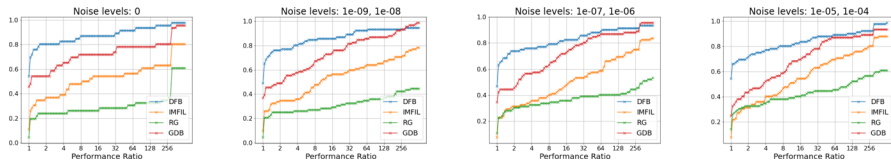


Fig. 8 Performance profiles of finite-difference-based methods on minimizing $\mathcal{C}^{1,1}$ functions

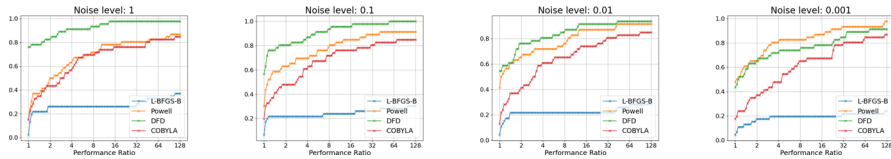


Fig. 9 Performance profiles of derivative-free methods on $\mathcal{C}^{1,1}$ functions with large, known noise levels

6.2 SciPy production-ready algorithms for functions with large noise

This subsection contains some illustrations of the performance of DFD (Algorithm 4) on the same subset of CUTEst problems [27] with the details given in Table 2. To demonstrate the efficiency of DFD in handling large noise, we inject the uniformly distributed stochastic noise into the tested problems as in the previous experiments with the high levels of noise $\xi_f \in \{1, 10^{-1}, 10^{-2}, 10^{-3}\}$. In this experiment, the performance of DFD is compared with *efficient production-ready codes* from the well-known SciPy library [66] of Python; namely, L-BFGS-B, Powell, and COBYLA algorithms. To the best of our knowledge, these methods are among the most popular, efficient, and state-of-the-art derivative-free methods for smooth functions. Although the Nelder-Mead method is also presented in the SciPy library, we do not consider it here due to its poor performance on smooth functions, since it does not take smooth structures into account in the algorithmic design.

All the algorithms are executed until they reach the maximum number of function evaluations of $200 \cdot n$. The setting of DFD is similar to the one given in Appendix A, while the settings of L-BFGS-B, Powell, and COBYLA algorithms are chosen to be standard without any modifications.

The illustration of the results is similar to the one mentioned in Subsection 6.1.2 and is presented in Figure 9. While we found that the Powell and COBYLA algorithms usually work well for the smallest noise $\xi_f = 10^{-3}$, our DFD method exhibits better results when the noise is larger, i.e., $\xi_f \geq 10^{-2}$. For this reason, we illustrate in Figure 10 below the results for few representative problems in 100-dimensional spaces with the noise levels 1 and 10^{-3} . Since the L-BFGS-B method does not achieve a comparable performance with other methods due to the large noise, we do not plot the results obtained by L-BFGS-B.

Note that the noise level is required for the implementation of DFD, which may not always be available in practice. However, this issue can be addressed if the noise is independent and identically distributed since a simple noise estimation procedure can be further applied. This assumption is also employed in the recent publications [63, Section 4], [6, Section 5]. In the following examples, we suppose that the noise level

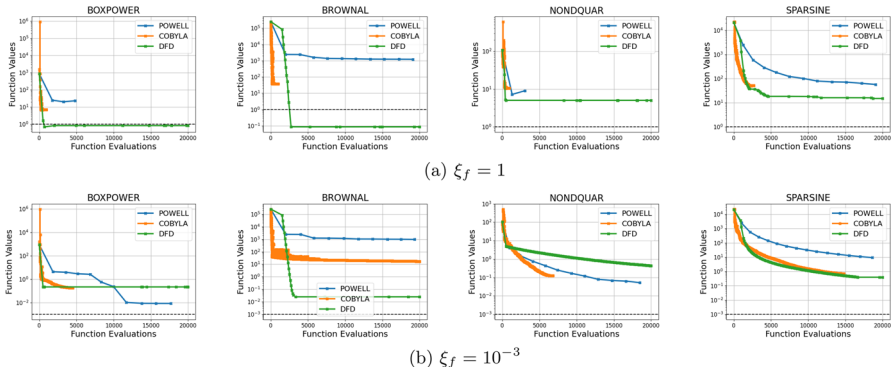


Fig. 10 Comparison of DFD with Powell and COBYLA algorithms from SciPy library. The exact function values against the function evaluations are presented. The dashed black line shows the noise level ξ_f of the function

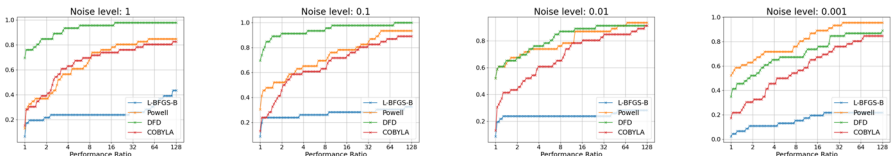


Fig. 11 Performance profiles of derivative-free methods on $C^{1,1}$ functions with large, unknown noise

is unknown for DFD and estimate the local noise level at an arbitrary point $x \in \mathbb{R}^n$ but using it as a global noise level. Given a sampling radius Δ and a sampling number $m \in \mathbb{N}$, let $u^i, i = 1, \dots, m$, be uniformly sampled in $\mathbb{B}(x, \Delta)$ and then define the computation noise level as

$$\varepsilon_f := \max_{i=1, \dots, m} \left\{ f(u^i) - \frac{1}{m} \sum_{j=1}^m f(u^j) \right\}.$$

In the following experiment, we choose $\Delta = 10^{-15}$ and $m = 2n$, where n is the dimension of the objective function. The results show that DFD still significantly outperforms all other methods when the noise is large and has comparable performance to the Powell method when the noise level is smaller.

6.3 Acceleration techniques for $\mathcal{C}_L^{1,1}$ functions

In this section, we present numerical enhancements to Algorithm 2 by incorporating acceleration techniques. First, we consider the Polyak heavy-ball method, also known as the Polyak momentum method [56], which has been well recognized in the literature on smooth convex and nonconvex optimization. It has been justified to achieve a faster convergence rate compared to standard gradient descent for minimizing $\mathcal{C}_L^{1,1}$ strongly convex functions in [56]. The convergence analysis for nonconvex smooth objective

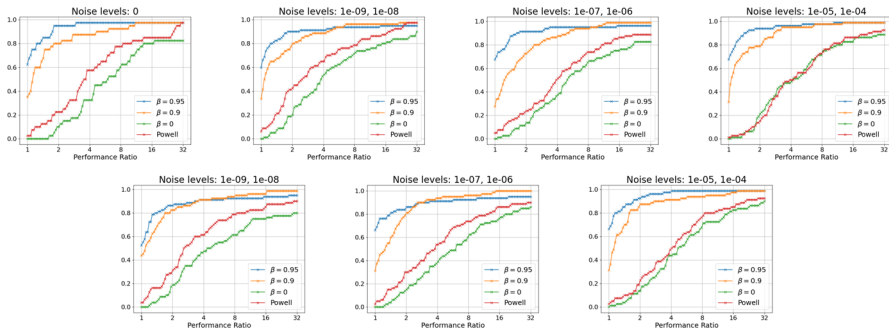


Fig. 12 Performance profiles of DFC-HB and Powell methods on minimizing $C_L^{1,1}$ functions with i.i.d. uniform noise (upper) and non-i.i.d. noise (lower)

functions is also discussed in [68] with further extensions to nonsmooth optimization methods given in [67]. In what follows, we numerically examine the integration of this method into DFC via Algorithm 2 for solving convex and nonconvex noisy smooth problems in the derivative-free setting, while deferring a rigorous theoretical analysis in this context. To this end, Algorithm 2 is modified in the following way, where "HB" is the abbreviation of the "heavy ball".

Algorithm 5 (DFC-HB).

Step 0 Select some $x^1 \in \mathbb{R}^n$, $\delta_1 > 0$, $L_1 > 0$, $\theta \in (0, 1)$, $\beta > 0$, and $\eta > 1$.

Step 1 (approximate gradient). Find g^k and the smallest nonnegative integer i_k such that

$$g^k = \tilde{G}(x^k, \theta^{i_k} \delta_k) \quad \text{and} \quad \|g^k\| > 2L_k \sqrt{n} \theta^{i_k} \delta_k. \quad (6.1)$$

Then set $\delta_{k+1} := \theta^{i_k} \delta_k$.

Step 2 (update). If $\phi\left(x^k - \frac{1}{L_k} g^k\right) \leq \phi(x^k) - \frac{1}{24L_k} \|g^k\|^2$, then $x^{k+1} := x^k + \beta(x^k - x^{k+1}) - \frac{1}{L_k} g^k$ and $L_{k+1} := L_k$. Otherwise, set $x^{k+1} := x^k$ and $L_{k+1} := \eta L_k$.

Note that DFC-HB is a generalization of DFC with $\beta = 0$ reducing it to the standard DFC without momentum. Using the same experimental settings as those in Subsections 6.1.1 and 6.2, we present the results for DFC-HB in comparison with the standard DFC and the Powell method below. The results in Figure 12 show that, while the standard version of DFC does not perform well compared to the Powell method, incorporating momentum significantly improves its performance. As a result, DFC with momentum ($\beta = 0.9$ and $\beta = 0.95$) outperforms the Powell method in this numerical experiment across different noise levels and types. Note that the values $\beta = 0.9$ and $\beta = 0.95$ are not derived from an extensive grid search, but are standard

momentum selections widely used in machine learning tasks, e.g., [43, Section 5] and the references therein. The data profiles for the methods are presented in Appendix B.

We now turn our attention to *quasi-Newton methods* [53, Chapter 6], a very successful technique for nonlinear continuous optimization that approximates the Hessian matrix to the navigate curvature of the objective function without much of the computational cost of exact Newton methods. The general framework of quasi-Newton methods when incorporated with DFC is as follows, where the abbreviation “QN” stands for quasi-Newton.

Algorithm 6 (DFC-QN).

Step 0. Select some $x^1 \in \mathbb{R}^n$, $\delta_1 > 0$, $L_1 > 0$, $\theta \in (0, 1)$, and $\eta > 1$.

Step 1 (approximate gradient). Find g^k and the smallest nonnegative integer i_k such that

$$g^k = \tilde{G}(x^k, \theta^{i_k} \delta_k) \quad \text{and} \quad \|g^k\| > 2L_k \sqrt{n} \theta^{i_k} \delta_k. \quad (6.2)$$

Then set $\delta_{k+1} := \theta^{i_k} \delta_k$.

Step 2 (update). If $\phi\left(x^k - \frac{1}{L_k} g^k\right) \leq \phi(x^k) - \frac{1}{24L_k} \|g^k\|^2$, then find $d^k \in \mathbb{R}^n$ such that $H_k d^k = -g^k$, where H_k is a Hessian approximation, and $t_k > 0$ such that

$$t_k := \max \left\{ t \mid f(x^k + t d^k) \leq f(x^k) - \beta t \|d^k\|^2 \mid t = 1, \gamma, \gamma^2, \dots \right\},$$

and set $x^{k+1} := x^k + t_k d^k$ and $L_{k+1} := L_k$. Otherwise, $x^{k+1} := x^k$ and $L_{k+1} := \eta L_k$.

In this numerical experiment, we consider two most well-known types of quasi-Newton methods in the literature, which are BFGS (Broyden-Fletcher-Goldfarb-Shanno) and L-BFGS (limited-memory BFGS). The sequence $\{H_k\}_{k \in \mathbb{N}}$ using the BFGS updates can be computed as follows: starting from a positive definite matrix H_1 , for each $k \in \mathbb{N}$, we define vectors

$$s^k := x^{k+1} - x^k, \quad y^k := g^{k+1} - g^k,$$

where $x^k, x^{k+1}, g^k, g^{k+1}$ are taken from Algorithm 6. Then the updated matrix is defined as

$$H_{k+1} := \begin{cases} H_k + \frac{y_k y_k^T}{\langle y^k, s^k \rangle} - \frac{H_k s^k (H_k s^k)^T}{\langle H_k s^k, s^k \rangle}, & \text{if } \langle s^k, y^k \rangle > 0, \\ H_k, & \text{otherwise.} \end{cases}$$

Note that no matrix inversion is needed to compute d^k from equation $H_k d^k = -g^k$ in practice, since it is possible to construct the sequence $\{H_k^{-1}\}_{k \in \mathbb{N}}$ iteratively; see [53, Equation (6.17)].

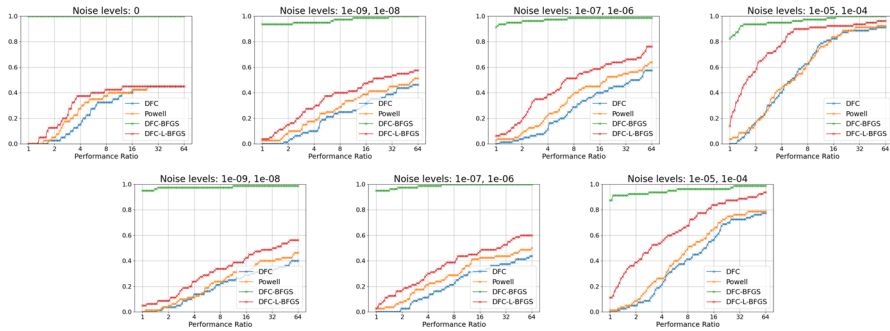


Fig. 13 Performance profiles of DFC-QN and Powell methods on minimizing $C_L^{1,1}$ functions with uniform i.i.d. noise (upper) and non-i.i.d. noise (lower)

When dealing with a large number of variables, L-BFGS updates are also useful since they reduce the cost of storing and updating approximations of the Hessian matrix. This is now considered as probably the most widely used method of this class, which was first introduced in [48]. It is based on the BFGS update employed at iteration k only to the most recent $\min\{m, k\}$ pairs (here m is a parameter, usually chosen from $\{3, \dots, 20\}$) to compute a descent direction.

Using the same settings as in Subsection 6.1.1 and Subsection 6.2, the results of DFC-BFGS and DFC-L-BFGS, in comparison with the standard DFC and the state-of-the-art Powell method from the Scipy library, are presented in Figure 13 below. Similarly to the results above for momentum techniques, DFC while incorporating with quasi-Newton steps, also outperforms the Powell method in most scenarios. The data profiles for the methods are presented in Appendix B.

7 Concluding remarks

This paper addresses derivative-free optimization problems with smooth and not necessarily convex objectives. A general derivative-free optimization method with a constant stepsize (DFC) is proposed to deal with $C_L^{1,1}$ problems. This novel method is shown to achieve the fundamental convergence properties of the standard gradient descent in the noiseless case and reach a near-stationary point in the noisy case without demanding any noise level information. Constructive estimates of the number of required iterations and function evaluations are established in the paper.

To deal with $C_L^{1,1}$ problems, a general derivative-free optimization method with backtracking stepsize (DFB) is proposed. The analysis of DFB in the noiseless case recovers convergence properties of the standard gradient descent method with a backtracking stepsize. To handle $C_L^{1,1}$ problems with large noise, a derivative-free optimization method with dynamic step linesearch (DFD) is proposed. It is revealed that DFD offers greater robustness than other finite-difference-based schemes to solve $C_L^{1,1}$ problems with complex structure. The conducted analysis shows that under certain conditions, DFD reaches a near-stationary point after a finite number of iterations.

Numerical results demonstrate that DFC and DFB achieve higher efficiency and robustness in comparison with other well-known finite-difference-based schemes in solving noiseless problems and problems with small noise. Moreover, DFD provide favorable results compared to some production-ready codes from SciPy library when the noise is large.

Our future research includes convergence analysis of the newly developed algorithms coupled with quasi-Newton methods for noisy smooth functions together with the appropriate accelerations, specifically the ones presented in the numerical experiments in Section 6. Moreover, we intend to establish efficient conditions to ensure local and global convergence to local minimizers of iterative sequences generated by derivative-free methods for problems of nonsmooth unconstrained and constrained optimization. Further practical conditions with rigorous theoretical guarantees to satisfy the requirements of Theorem 4.10 will also be considered in our future research.

Acknowledgements We sincerely thank two anonymous referees for their valuable time to referee the paper and insightful feedback that allowed us to improve the original presentation. Our gratitude also goes to Katya Scheinberg for many fruitful discussions and remarks; in particular, for her suggestions on acceleration techniques, which helped us to significantly enhance the quality and practicality of our paper.

Declarations

Conflict of interest The authors declare that they have no conflict of interest.

A Numerical results on bivariate functions

In this appendix, we present additional results for the experiment conducted in Remark 5.7. The setup for the experiment is the following:

- DFD: The parameters are chosen as $\eta = 2$, $L_1 = 1$.
- IMFIL: The setting of IMFIL in this experiment follows the original development at [26, Page 279], with $\bar{\alpha} = 10^{-10}$; $\beta = 0.1$; $\gamma = 0.5$ and $h_k = 2^{1-k}$.
- RG: The parameters of RG in this experiment also obey the equations (55) and (58) in the original paper [52], i.e., $h = \frac{1}{4(n+4)L}$ and $\mu = \frac{5}{3(n+4)}\sqrt{\frac{\varepsilon}{2L}}$. Since the function in question does not have a globally Lipschitz continuous gradient, we tune the Lipschitz constant L by grid search on the set $\{0.1, 1, 10\}$ and choose the best one corresponding to the smallest function value at the last iterate.
- DF-backtracking: The parameters are chosen as $\eta = 2$ and $L_1 = 1$ similarly to DFD.
- GDD (Ada): The code for the adaptive finite difference interval estimation is given in [63, Algorithm 2.1]. The parameters for dynamic step linesearch are similar to DFD.
- L-BFGS (Ada): The L-BFGS code¹ is provided is taken from [64], while the code for the adaptive finite difference interval estimation is provided by [63, Algorithm 2.1].

¹ <https://github.com/hjmshi/noise-tolerant-bfgs>

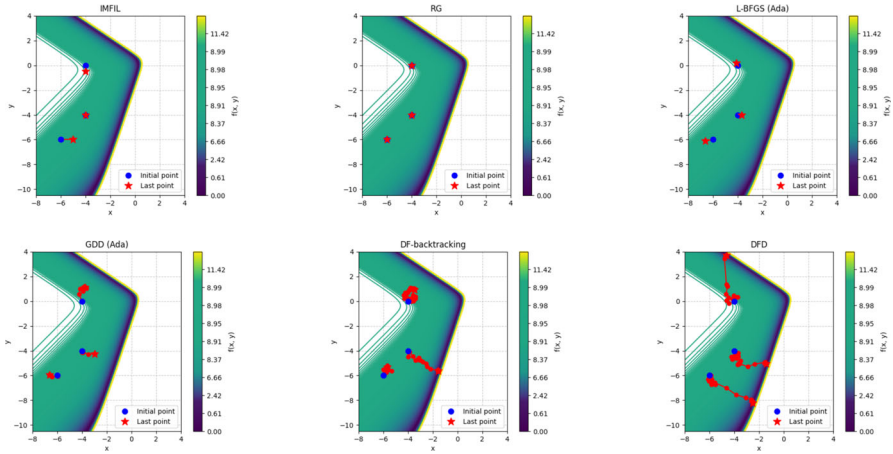


Fig. 14 Finite-difference-based methods on minimizing a bivariate $C^{1,1}$ function (noise level 1)

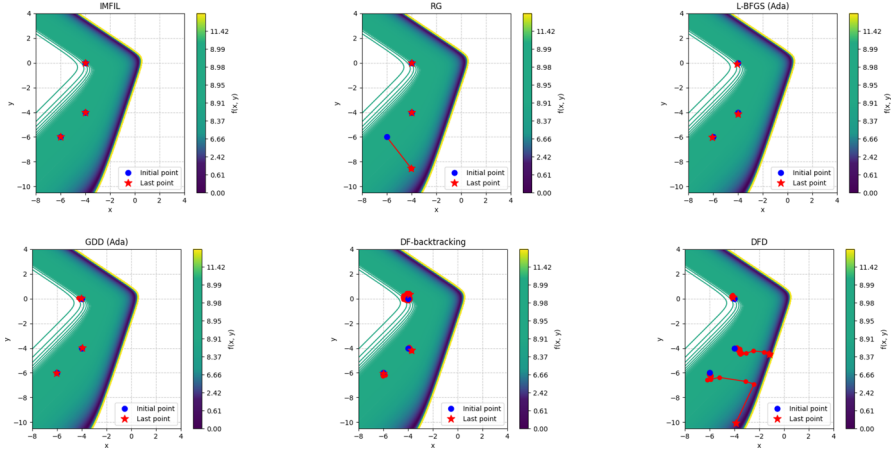


Fig. 15 Finite-difference-based methods on minimizing a bivariate $C^{1,1}$ function (noise level 10^{-1})

All the algorithms are executed for 200 function evaluations with the three different initial points $(-4, 0)$, $(-4, -4)$, $(-6, 0)$. We also choose the noise levels from the set $\{1, 0.1, 0.01, 0.001\}$. Since the result with a noise level of 0.01 is already presented in Remark 5.7, we do not represent it here. In addition, while conducting the experiments, due to the randomness of the noisy objective function, there are some cases where the iterative sequence generated by the RG method explodes to extremely large numbers (around 10^{26}) and does not find the minimum region properly. For this reason, we exclusively plot points generated by the methods within a ball centered at the origin with the radius 20. It can be seen that our DFD is stable with respect to different levels of noise, and fails only in one over nine cases when the noise is 0.1 and the initial point is $(-4, 0)$ (Figures 14, 15, 16).

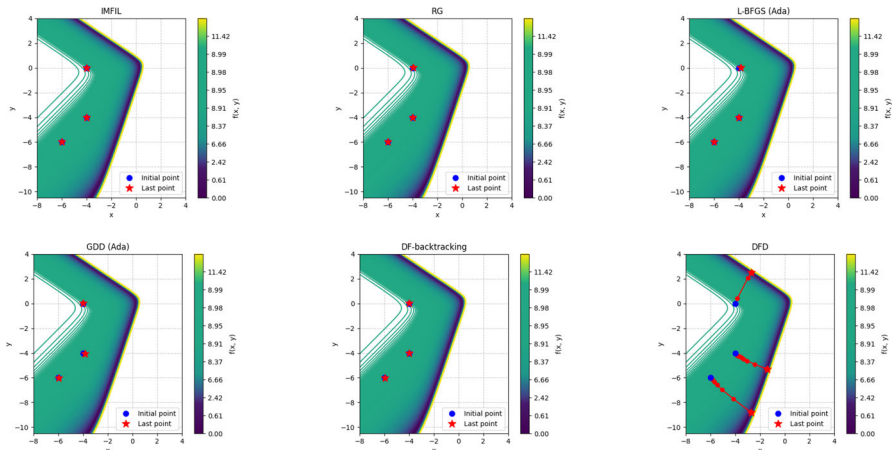


Fig. 16 Finite-difference-based methods on minimizing a bivariate $C^{1,1}$ function (noise level 10^{-3})

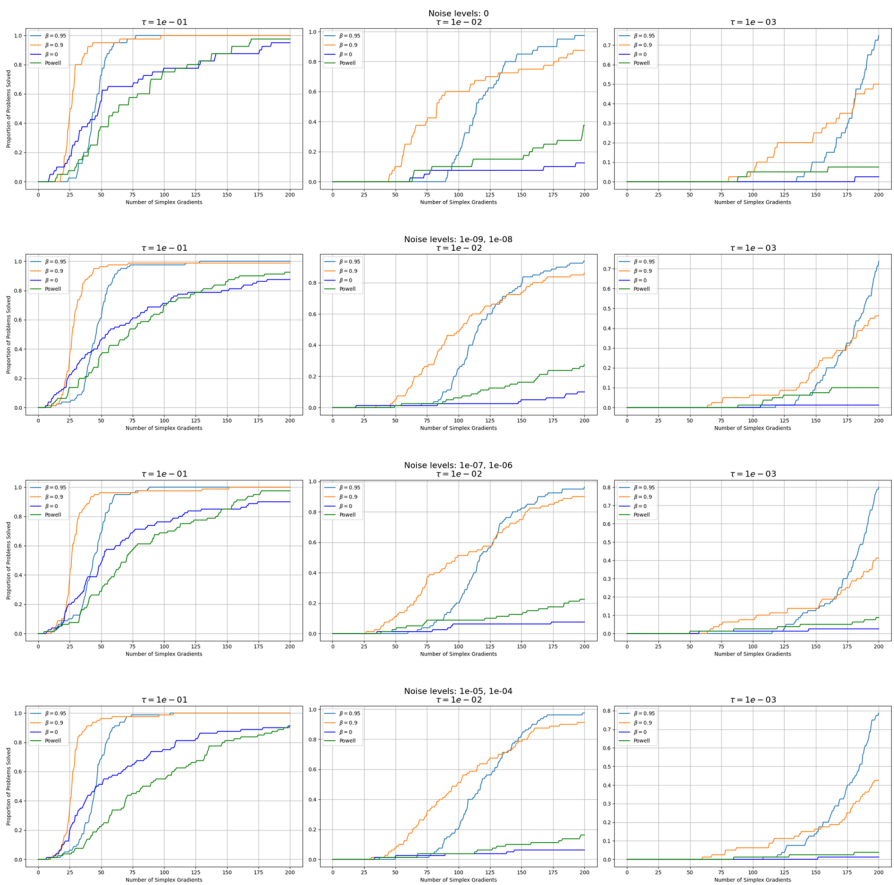


Fig. 17 Data Profiles of DFC-HB and Powell methods on minimizing $C_L^{1,1}$ functions with uniform i.i.d. noise

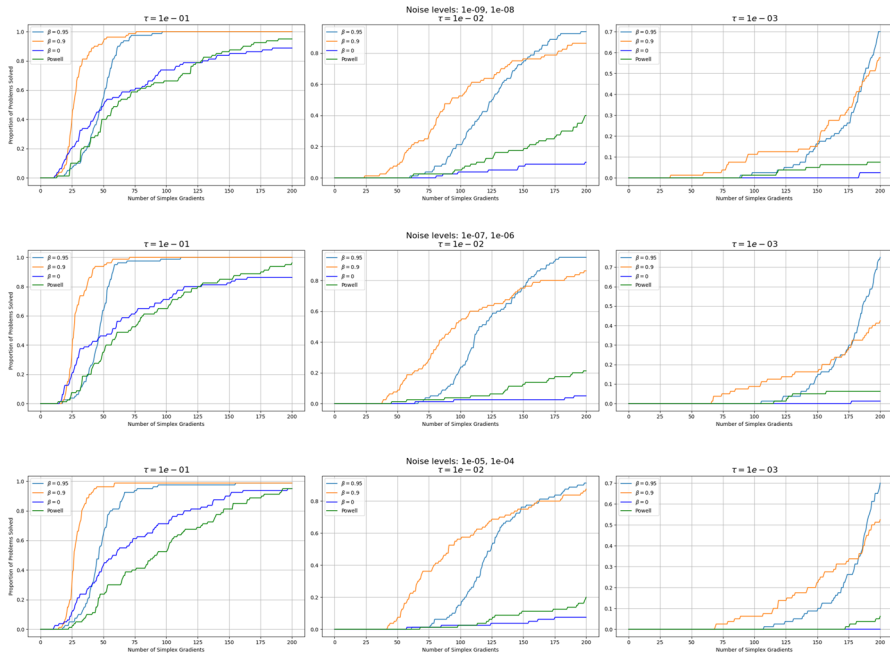


Fig. 18 Data Profiles of DFC-HB and Powell methods on minimizing $C_L^{1,1}$ functions with non-i.i.d. noise

B Data profiles for methods in section 6.3

In this appendix, we present comparisons between the methods in Section 6.3 using data profiles [50], where a solver s is said to solve a problem p with accuracy τ if it reaches a point x_s^* such that

$$f_p(x_s^1) - f_p(x_s^*) \geq (1 - \tau)(f_p(x_s^1) - f_p^*),$$

where f_p is the objective function of problem p , x_s^1 is the initial point, and f_p^* is the smallest function value obtained by all solvers for problem p . The “number of simplex gradients” for each solver is defined as the number of $(n + 1)$ -bundles of gradients used. The algorithm settings are similar to those in Section 6.3, except that the maximum number of function evaluations is set to $200(n + 1)$ instead of $200n$. The results also demonstrate that the acceleration techniques significantly improve the performance of the DFC method and, in most cases, outperform the state-of-the-art Powell method from the SciPy library in terms of computational budget (Figures 17, 18, 19, 20).

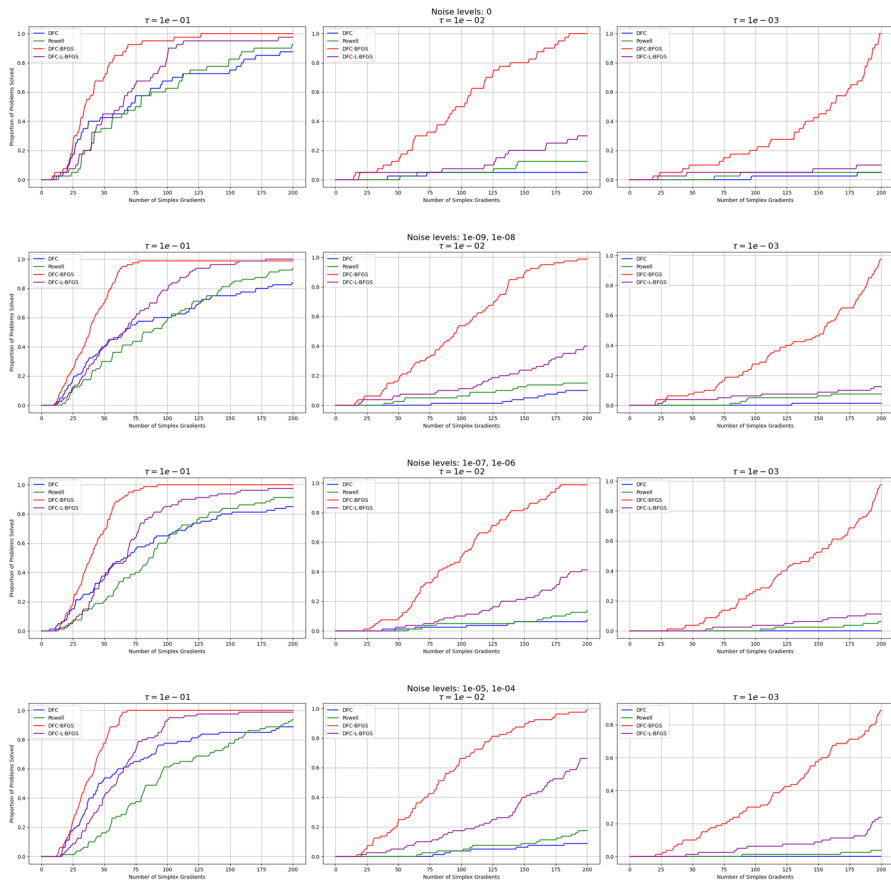


Fig. 19 Data Profiles of DFC-QN and Powell methods on minimizing $C_L^{1,1}$ functions with uniform i.i.d. noise

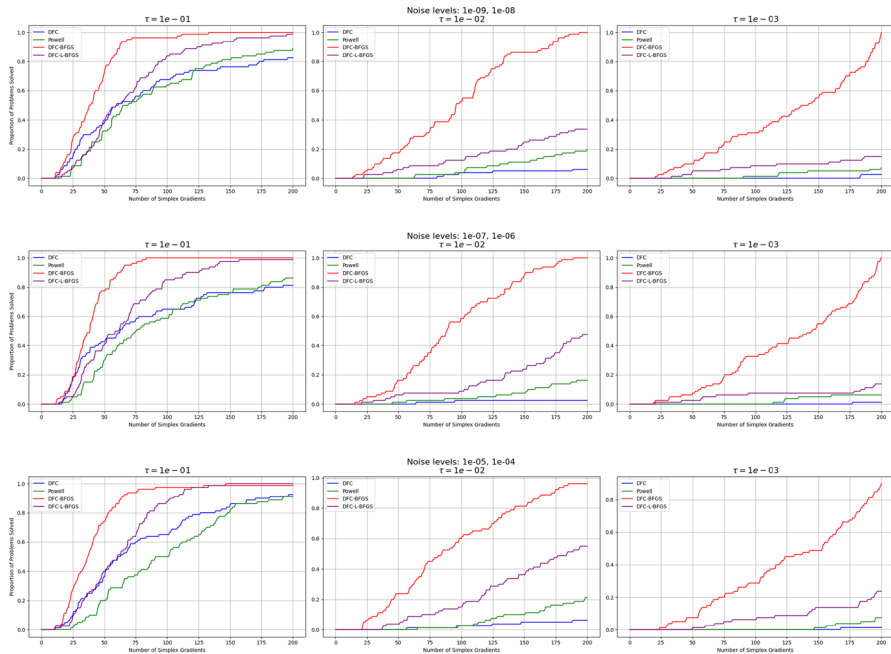


Fig. 20 Data Profiles of DFC-QN and Powell methods on minimizing $C_L^{1,1}$ functions with non-i.i.d. noise

References

1. Absil, P.-A., Mahony, R., Andrews, B.: Convergence of the iterates of descent methods for analytic cost functions. *SIAM J. Optim.* **16**, 531–547 (2005)
2. Addis, A., Cassioli, A., Locatelli, M., Schoen, F.: A global optimization method for the design of space trajectories. *Comput. Optim. Appl.* **48**, 635–652 (2011)
3. Attouch, H., Bolte, J., Redont, P., Soubeyran, A.: Proximal alternating minimization and projection methods for nonconvex problems. An approach based on the Kurdyka-Łojasiewicz property. *Math. Oper. Res.* **35**, 438–457 (2010)
4. Attouch, H., Bolte, J., Svaiter, B.F.: Convergence of descent methods for semi-algebraic and tame problems: Proximal algorithms, forward-backward splitting, and regularized Gauss-Seidel methods. *Math. Program.* **137**, 91–129 (2013)
5. Audet, C., Hare, W.: *Derivative-Free and Blackbox Optimization*. Springer, Cham, Switzerland (2017)
6. Berahas, A.S., Byrd, R.H., Nocedal, J.: Derivative-free optimization of noisy functions via quasi-Newton methods. *SIAM J. Optim.* **29**, 965–993 (2019)
7. Berahas, A.S., Cao, L., Choromanski, K., Scheinberg, K.: A theoretical and empirical comparison of gradient approximations in derivative-free optimization. *Found. Comput. Math.* **22**, 507–560 (2022)
8. Berahas, A.S., Cao, L., Scheinberg, K.: Global convergence rate analysis of a generic linesearch algorithm with noise. *SIAM J. Optim.* **31**, 1489–1518 (2021)
9. Bertsekas, D.P.: *Nonlinear Programming*, 3rd edn. Athena Scientific, Belmont, MA (2016)
10. Bertsekas, D.P., Tsitsiklis, J.N.: Gradient convergence in gradient methods with errors. *SIAM J. Optim.* **10**, 627–642 (2000)
11. Bolte, J., Pauwels, E.: Conservative set-valued fields, automatic differentiation, stochastic gradient methods and deep learning. *Math. Program.* **188**, 19–51 (2021)
12. Bellavia, S., Gurioli, G., Morini, B., Toint, P.L.: Adaptive regularization algorithms with inexact evaluations for nonconvex optimization. *SIAM J. Optim.* **29**, 2881–2915 (2019)
13. Cartis, C., Gould, N.I.M., Toint, P.: On the oracle complexity of first-order and derivative-free algorithms for smooth nonconvex minimization. *SIAM J. Optim.* **22**, 66–86 (2012)

14. Cartis, C., Scheinberg, K.: Global convergence rate analysis of unconstrained optimization methods based on probabilistic models. *Math. Program.* **169**, 337–375 (2018)
15. Cartis, C., Gould, N.I.M., Toint, P.L.: *Evaluation Complexity of Algorithms for Nonconvex Optimization: Theory, Computation and Perspectives*, SIAM, Philadelphia, PA (2022)
16. Choi, T.D., Kelley, C.T.: Superlinear convergence and implicit filtering. *SIAM J. Optim.* **10**, 1149–1162 (2000)
17. Conn, A.R., Scheinberg, K., Vicente, L.N.: *Introduction to Derivative-Free Optimization*. SIAM, Philadelphia, PA (2009)
18. Conn, A.R., Scheinberg, K., Vicente, L.N.: Global convergence of general derivative-free trust-region algorithms to first- and second-order critical points. *SIAM J. Optim.* **20**, 387–415 (2009)
19. Dolan, E., Moré, J.: Benchmarking optimization software with performance profiles. *Math. Program.* **91**, 201–213 (2022)
20. Doikov, N., Grapiglia, G. N.: First and zeroth-order implementations of the regularized Newton method with lazy approximated Hessians, <https://doi.org/10.48550/arXiv.2309.02412>
21. Duchi, J.C., Jordan, M.I., Wainwright, M.J., Wibisono, A.: Optimal rates for zero-order convex optimization: The power of two function evaluations. *IEEE Trans. Inf. Theory* **61**, 2788–2806 (2015)
22. Facchinei, F., Pang, J.-S.: *Finite-Dimensional Variational Inequalities and Complementarity Problems*, vol. II. Springer, New York (2003)
23. Fowkes, J., Roberts, L., Burmen, A.: PyCUTEst: An open source Python package of optimization test problems. *J. Open Source Softw.* **7**, 4377 (2022)
24. Fridovich-Keil, S., Recht, B.: Choosing the stepsize: Intuitive linesearch algorithms with efficient convergence, The 11th Workshop on Optimization for Machine Learning (2019)
25. Gill, P.E., Murray, W., Saunders, M.A., Wright, M.H.: Computing forward-difference intervals for numerical optimization. *SIAM J. Sci. Comput.* **4**, 310–321 (1983)
26. Gilmore, P., Kelley, C.T.: An implicit filtering algorithm for optimization of functions with many local minima. *SIAM J. Optim.* **5**, 269–285 (1995)
27. Gould, N.I., Orban, D., Toint, P.L.: CUTEst: A constrained and unconstrained testing environment with safe threads for mathematical optimization. *Comput. Optim. Appl.* **60**, 545–557 (2015)
28. Gasnikov, A., Dvinskikh, D., Dvurechensky, P., Gorbunov, E., Beznosikov, A., Lobanov, A.: Randomized gradient-free methods in convex optimization, arXiv preprint [arXiv:2211.13566](https://arxiv.org/abs/2211.13566), (2022)
29. Gorbunov, E., Dvurechensky, P., Gasnikov, A.: An accelerated method for derivative-free smooth stochastic convex optimization. *SIAM J. Optim.* **32**, 1210–1238 (2022)
30. Grapiglia, G.N.: Worst-case evaluation complexity of a derivative-free quadratic regularization method. *Optim. Lett.* **18**, 195–213 (2024)
31. Gray, G.A., Kolda, T.G.: Algorithm 856: Appspack 4.0: Asynchronous parallel pattern search for derivative-free optimization. *ACM Trans. Math. Softw.* **32**, 485–507 (2006)
32. Hare, W., Lucet, Y.: Derivative-free optimization via proximal point methods. *J. Optim. Theory Appl.* **160**, 204–220 (2014)
33. Hare, W., Nutini, J.: A derivative-free approximate gradient sampling algorithm for finite minimax problems. *Comput. Optim. Appl.* **56**, 1–38 (2013)
34. Hooke, R., Jeeves, T.A.: Direct search solution of numerical and statistical problems. *J. ACM* **8**, 212–229 (1961)
35. Izmailov, A.F., Solodov, M.V.: *Newton-Type Methods for Optimization and Variational Problems*. Springer, New York (2014)
36. Josz, C.: Global convergence of the gradient method for functions definable in o-minimal structures. *Math. Program.* **202**, 355–383 (2023)
37. Josz, C., Lai, L., Li, X.: Convergence of the momentum method for semi-algebraic functions with locally Lipschitz gradients. *SIAM J. Optim.* **33**, 2988–3011 (2023)
38. Kelley, C.T.: *Iterative Methods for Optimization*. SIAM, Philadelphia, PA (1999)
39. Kelley, C. T.: *Implicit Filtering and Nonlinear Least Squares Problems*, The International Federation for Information Processing, 71–90 (2003)
40. Khanh, P.D., Mordukhovich, B.S., Tran, D.B.: Inexact reduced gradient methods in smooth nonconvex optimization. *J. Optim. Theory Appl.* **203**, 2138–2178 (2024)
41. Khanh, P.D., Mordukhovich, B.S., Tran, D.B.: A new inexact gradient descent method with applications to nonsmooth convex optimization. *Optim. Methods Softw.* (2024). <https://doi.org/10.1080/10556788.2024.2322700>

42. Khanh, P.D., Mordukhovich, B.S., Phat, V.T., Tran, D.B.: Inexact proximal methods for weakly convex functions. *J. Glob. Optim.* **91**, 611–646 (2025)
43. Khanh, P.D., Luong, H.-C., Mordukhovich, B.S., Tran, D.B.: Fundamental convergence analysis of sharpness-aware minimization. *Adv. Neural Inf. Process. Syst.* (2024)
44. Kurdyka, K.: On gradients of functions definable in o-minimal structures. *Ann. Inst. Fourier* **48**, 769–783 (1998)
45. Lagarias, J.C., Reeds, J.A., Wright, M.H., Wright, P.E.: Convergence properties of the Nelder-Mead simplex method in low dimensions. *SIAM J. Optim.* **9**, 112–147 (1998)
46. Lee, J.D., Panageas, I., Piliouras, G., Simchowitz, M., Jordan, M.I., Recht, B.: First-order methods almost always avoid strict saddle points. *Math. Program.* **176**, 311–337 (2019)
47. Łojasiewicz, S.: *Ensembles Semi-Analytiques*. Institut des Hautes Etudes Scientifiques, Bures-sur-Yvette, France (1965)
48. Liu, D.C., Nocedal, J.: On the limited memory BFGS method for large scale optimization. *Math. Program.* **45**, 503–528 (1989)
49. Moré, J.J., Wild, S.M.: Estimating derivatives of noisy simulations. *ACM Trans. Math. Softw.* **38**, 1–21 (2012)
50. Moré, J.J., Wild, S.M.: Benchmarking derivative-free optimization algorithms. *SIAM J. Optim.* **20**, 172–191 (2009)
51. Nelder, J.A., Mead, R.: A simplex method for function minimization. *Comput. J.* **7**, 308–313 (1965)
52. Nesterov, Y., Spokoiny, V.: Random gradient-free minimization of convex functions. *Found. Comput. Math.* **17**, 527–566 (2017)
53. Nocedal, J., Wright, S. J.: *Numerical Optimization*, 2nd edition. New York, (2006)
54. Ostrowski, A.: *Solution of Equations and Systems of Equations*, 2nd edn. Academic Press, New York (1966)
55. Polyak, B.T.: Gradient methods for the minimization of functionals. *USSR Comput. Math. Math. Phys.* **3**, 864–878 (1963)
56. Polyak, B.T.: Some methods of speeding up the convergence of iteration methods. *USSR Comput. Math. Math. Phys.* **3**, 1–17 (1964)
57. Powell, M.J.D.: An efficient method for finding the minimum of a function of several variables without calculating derivatives. *Comput. J.* **7**, 155–162 (1964)
58. Powell, M. J. D.: A direct search optimization method that models the objective and constraint functions by linear interpolation, in *Advances in Optimization and Numerical Analysis*, pp. 51–67, Springer, New York, (1994)
59. Themelis, A., Stella, L., Patrinos, P.: Forward–backward quasi-Newton methods for nonsmooth optimization problems. *Comput. Optim. Appl.* **67**, 443–487 (2017)
60. Truong, T.T., Nguyen, H.-T.: Backtracking gradient descent method and some applications in large scale optimisation. Part 2: Algorithms and experiments. *Appl. Math. Optim.* **84**, 2557–2586 (2021)
61. Scheinberg, K.: Finite difference gradient approximation: To randomize or not? *INFORMS J. Comput.* **34**, 2384–2388 (2022)
62. Shibaev, I., Dvurechensky, P., Gasnikov, A.: Zeroth-order methods for noisy Hölder-gradient functions. *Optim. Lett.* **16**, 2123–2143 (2022)
63. Shi, H.M., Xie, Y., Xuan, M.Q., Nocedal, J.: Adaptive finite-difference interval estimation for noisy derivative-free optimization. *SIAM J. Sci. Comput.* **44**, 2302–2321 (2022)
64. Shi, H.M., Xie, Y., Byrd, R., Nocedal, J.: A noise-tolerant quasi-Newton algorithm for unconstrained optimization. *SIAM J. Optim.* **32**, 29–55 (2022)
65. Shi, H.M., Xuan, M.Q., Oztoprak, F., Nocedal, J.: On the numerical performance of finite-difference-based methods for derivative-free optimization. *Optim. Methods Softw.* **38**, 289–311 (2023)
66. Virtanen, P., Gommers, R., Oliphant, T.E., et al.: *SciPy 1.0: Fundamental algorithms for scientific computing in Python*. *Nat. Methods* **17**, 261–272 (2020)
67. Wen, B., Chen, X., Pong, T.K.: Linear convergence of proximal gradient algorithm with extrapolation for a class of nonconvex nonsmooth minimization problems. *SIAM J. Optim.* **27**, 124–145 (2017)
68. Zavriev, S., Kostyuk, F.: Heavy-ball method in nonconvex optimization problems. *Comput. Math. Model.* **4**, 336–341 (1993)

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.