

On Approximations with Finite Precision in Bundle Methods for Nonsmooth Optimization¹

M. V. SOLODOV²

Communicated by P. Tseng

Abstract. We consider the proximal form of a bundle algorithm for minimizing a nonsmooth convex function, assuming that the function and subgradient values are evaluated approximately. We show how these approximations should be controlled in order to satisfy the desired optimality tolerance. For example, this is relevant in the context of Lagrangian relaxation, where obtaining exact information about the function and subgradient values involves solving exactly a certain optimization problem, which can be relatively costly (and as we show, in any case unnecessary). We show that approximation with some finite precision is sufficient in this setting and give an explicit characterization of this precision. Alternatively, our result can be viewed as a stability analysis of standard proximal bundle methods, as it answers the following question: for a given approximation error, what kind of approximate solution can be obtained and how does it depend on the magnitude of the perturbation?

Key Words. Nonsmooth optimization, convex optimization, bundle methods, stability analysis, perturbations.

1. Introduction

We consider the problem

$$\min_x \{f(x) | x \in \mathbb{R}^n\}, \quad (1)$$

where $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is a convex (in general, nondifferentiable) function. Bundle methods [see Ref. 1, Ref. 2 (Chapters 14–15), Ref. 3 (Chapter 7)] are perhaps

¹This work was partially supported by CNPq Grant 300734/95-6, PRONEX Optimization, and FAPERJ.

²Associate Researcher, Instituto de Matemática Pura e Aplicada, Rio de Janeiro, Brazil.

the most practical and efficient algorithms for solving problems of this class. We are interested in the case where, at each given point, we have available only an approximate value of f and an approximate value of one of its subgradients. For example, this can be considered a typical situation in the context of Lagrangian relaxation (Refs. 3–6), where one would need an exact solution of a certain optimization problem to obtain exact information about f and its subgradient. To make this motivation a little more specific, consider the (primal) problem

$$\max_{\xi} \{q(\xi) \mid \xi \in P, h(\xi) = 0\}, \quad (2)$$

where P is a compact subset of \mathbb{R}^m and $q: \mathbb{R}^m \rightarrow \mathbb{R}$, $h: \mathbb{R}^m \rightarrow \mathbb{R}^n$. Lagrangian relaxation of the equality constraints in this problem leads to problem (1), where

$$f(x) = \max_{\xi} \{q(\xi) + \langle x, h(\xi) \rangle \mid \xi \in P\} \quad (3)$$

is the dual function. Trying to solve problem (2) by means of solving its dual (1) makes sense in many situations. One example is when some separable structure is present in (2), so that the maximization problem in (3) is decomposable [e.g. see Ref. 3 (Chapter 8) and Ref. 6 (Chapter 6)]. The approach of Lagrangian relaxation gives rise to a convex nonsmooth problem (1) independently of any assumptions about (2). Of course, without further assumptions, there can be a duality gap between the primal and the dual problems. Another important issue is recovering primal solutions after solving the dual. Here, we shall not deal with these questions. At the very least, the optimal value of (1) provides always an upper bound for the optimal value of (2), which can be useful already by itself in many applications.

As is well known and easy to see, if for a given $x \in \mathbb{R}^n$ some $\xi(x)$ is a solution of the optimization problem in (3), then we have not only the dual function value

$$f(x) = q(\xi(x)) + \langle x, h(\xi(x)) \rangle,$$

but also a subgradient

$$h(\xi(x)) \in \partial f(x).$$

The point relevant for the present paper is that evaluating the function value $f(x)$ and a subgradient $y \in \partial f(x)$ requires solving exactly the optimization problem in (3). In some cases, computing an exact solution can be cost prohibitive or at least cost inefficient. As we shall show, actually it is in some sense unnecessary. One can also argue that computing an exact solution is computationally unrealistic to start with, as any numerical method used for the maximization in (3) would return an approximate solution according to

some (finite) termination criterion. At any rate, it is probably desirable not to spend more time on solving (3) than it is really necessary. The question that we ask is therefore the following: Given the desired optimality tolerance $\Delta_{\text{opt}} > 0$ for problem (1), to which precision should we evaluate the values of f and its subgradients in order to guarantee that the bundle method using these inexact data would terminate with a point satisfying the given tolerance Δ_{opt} ?

We also pose the following related stability questions. Given some nonzero (and not necessarily tending to zero) approximation error, what are the convergence properties of the bundle method? What kind of approximate solutions of (1) can be obtained, and how do they depend on the approximation error? Stability analysis of computational algorithms with nonvanishing perturbations (for example, induced by noisy data) is of independent interest and importance (e.g. Refs. 7–8).

The problem setting that we consider here is essentially the same as in Ref. 9. Specifically, given $x \in \mathbb{R}^n$ and $\tilde{\epsilon} \geq 0$ (as a practical matter, actually $\tilde{\epsilon} > 0$), we assume that we can find some $\tilde{f} \in \mathbb{R}$ and $y \in \mathbb{R}^n$ such that

$$\begin{aligned} f(x) &\geq \tilde{f} \geq f(x) - \tilde{\epsilon}, \\ f(\zeta) &\geq l(\zeta) := \tilde{f} + \langle y, \zeta - x \rangle, \quad \forall \zeta \in \mathbb{R}^n. \end{aligned}$$

As discussed in Ref. 9, this setting is realistic in many applications. In the context of Lagrangian relaxation, these conditions amount to computing some $\xi \in P$ which is $\tilde{\epsilon}$ -optimal in the maximization problem in (3) and then taking

$$\tilde{f} = q(\xi) + \langle x, h(\xi) \rangle, \quad y = h(\xi).$$

As is seen easily, the stated conditions mean that

$$y \in \partial f_{\tilde{\epsilon}}(x),$$

where

$$\partial f_{\epsilon}(x) = \{y \in \mathbb{R}^n \mid f(\zeta) \geq f(x) + \langle y, \zeta - x \rangle - \epsilon\}$$

is the ϵ -subdifferential of f at x , with ∂f_0 being the exact subdifferential ∂f .

In Ref. 9, it was established that the proximal bundle method based on the inexact linearizations of f described above converges to a solution if the approximations of f and its subgradients are asymptotically exact, that is, if $\tilde{\epsilon} \rightarrow 0$ over the course of the iterations. This means that, for a given optimality tolerance $\Delta_{\text{opt}} > 0$, the method would terminate finitely with an acceptable approximate solution of (1) (also, at termination, $\tilde{\epsilon}$ would still be at some nonzero value). Note that this result says nothing about how to control the accuracy parameter $\tilde{\epsilon}$ and in particular about how small $\tilde{\epsilon}$ needs to be or

about how large it can be allowed to be, in order to ensure the finite termination. In the context of Lagrangian relaxation, it presumes that we can solve the optimization problem in (3) with an arbitrarily high precision. While this can be accepted as realistic in some cases, it may not be so in others. Also, it may be wasteful to approximate the objective and subgradient values [e.g., solve (3)] with higher precision than it is really necessary. Intuition suggests that it is reasonable to set this precision based on the required final optimality tolerance Δ_{opt} , especially at the later stages of computations. Making this relation explicit is the subject of this paper. In particular, we give an explicit expression for $\tilde{\epsilon}$ in terms of Δ_{opt} ; see (6).

Alternatively, our results provide stability analysis of the proximal bundle method, answering the following questions. Given some nonzero (and not tending to zero) approximation error $\tilde{\epsilon}$, what kind of approximate solutions of (1) can be obtained? How do they depend on the approximation error $\tilde{\epsilon}$? We comment on these issues further at the end of the paper.

As other work on bundle methods with approximate data, we mention Refs. 10–11. In Ref. 10, it is assumed that, for each $x \in \mathcal{R}^n$, the exact value $f(x)$ is available while the subgradient is computed approximately, $y \in \partial f_{\tilde{\epsilon}}(x)$. The attractive feature of the analysis in Ref. 10 is that, unlike in Ref. 9 and the present paper, no knowledge of ϵ is assumed. On the other hand, the exact value of f is needed, and so this setting is not suitable for some important applications of bundle methods and in particular for the Lagrangian relaxation. Indeed, as discussed above, in that case evaluating $f(x)$ and $y \in \partial f(x)$ is the same task, which is computing a solution $\xi(x)$ of the maximization problem in (3). In Ref. 11, an inexact proximal bundle method for the maximum eigenvalue function is developed. As in Ref. 9 and the present paper, both the function and subgradient values are approximated. Furthermore, certain rules to control the approximation are given. These rules are constructive and are preferable to merely saying that $\tilde{\epsilon} \rightarrow 0$. But as in Ref. 9, stability issues and finite (not tending to zero) precision of approximation are not considered in Ref. 11.

Our notation is standard. By $\langle \cdot, \cdot \rangle$ we denote the usual inner product in the given finite-dimensional space, with $\|\cdot\|$ being the associated norm. For an arbitrary set A , we denote by $|A|$ its cardinality.

2. Bundle Method with Inexact Data

Apart from being based on inexact data and a special stopping test, the algorithm that we consider here is a fairly standard proximal bundle method, with the size of the model controlled by the aggregation technique. Nevertheless, it is worth to explain our notation.

It is convenient for our purposes to keep separate track of the usual linear approximations to the function and of the approximations obtained after the aggregation. The former will be collected into the set B_k^c and denoted by $l_i(x)$, $i \leq k$. The latter will be collected into the set B_k^a and denoted by $l_i^a(x)$, $i \leq k$. Thus, B_k^c and B_k^a are subsets in the space of linear functions. The two sets I_k^c and I_k^a collect the iteration indices of the current members of B_k^c and B_k^a , respectively.

When the number of elements in $B_k^c \cup B_k^a$ reaches the prescribed upper bound B_{\max} , two or more of those elements are deleted from the bundle, being replaced by the so-called aggregate piece (Step 8 of Algorithm 2.1). This controls the complexity of the cutting-plane approximation of f [i.e., φ_k given by (8)] and thus keeps subproblem (7) in Step 2 of Algorithm 2.1 manageable. Subproblem (7) is solved via some quadratic programming method applied to its dual [see (14) and Lemma 3.1], whose dimensionality is precisely $|B_k^c| + |B_k^a|$. It makes sense to use specialized highly effective QP solvers developed for the structure in (14); e.g. see Ref. 12.

We do not give any specific rule for choosing the proximal parameter γ_k in Step 2 of Algorithm 2.1. This choice is important for the practical efficiency of any bundle method, but it is not central to the subject of this paper. For some rules to control γ_k , see Refs. 13, 14. Also, it is possible and can be desirable to use more sophisticated quadratic terms in (7). For some possibilities, see Refs. 15–20.

In Step 3 of Algorithm 2.1, the value of the predicted descent is computed and the solution of (7) is accepted as the next best iterate if the actual descent obtained with this point is at least a σ -fraction of the predicted (Step 6 of Algorithm 2.1). The indices of such descent steps are collected into the set K_d . The current best (approximate) value of f is recorded as \tilde{f}_k .

According to Proposition 3.1 below, the stopping criterion (10) in Step 4 of Algorithm 2.1 ensures that, if the method terminates with some iterate x^k , then we have

$$d^k \in \partial_{\epsilon_k} f(x^k) \quad (4)$$

such that

$$\Delta_k := (1/2\gamma_k) \|d^k\|^2 + \epsilon_k \leq \Delta_{\text{opt}}, \quad (5)$$

where d^k is computed using the solution of (14) [dual to (7), see Lemma 3.1]. Relations (4), (5) essentially constitute the usual stopping test in bundle methods; see e.g. Ref. 3 (Chapter 7). In Theorem 3.1 below, we show that our stopping test is guaranteed to be satisfied eventually if the approximations are controlled according to the following rule:

$$[(1 - \sigma)/2(2 - \sigma)]\Delta_{\text{opt}} > \limsup_k (\max\{\tilde{\epsilon}_i \mid i \in I_k^c\}), \quad (6)$$

where $\sigma \in (0, 1)$ is the relaxation parameter used in the descent test (Step 6 of Algorithm 2.1). In other words, if the approximation of the objective and subgradient values is (asymptotically) accurate enough in the sense of (6), then Algorithm 2.1 terminates finitely with a point x^k satisfying (4), (5). We emphasize that the required precision is finite ($\tilde{\epsilon}_k$ need not tend to zero) and that we show precisely how small/large it needs to be for computing the desired approximate solution of (1). At the end of the paper, our results are further translated into the language of stability analysis.

Algorithm 2.1. Choose some $\sigma \in (0, 1)$, $x^0 \in \mathfrak{R}^n$, and set $K_d, B_{-1}^c, B_0^a, I_0^a := \emptyset$. Choose an integer $B_{\max} \geq 2$. Compute some $\tilde{f}_0 \in \mathfrak{R}$ and $y^0 \in \mathfrak{R}^n$ such that $f(x^0) \geq \tilde{f}_0 \geq f(x^0) - \tilde{\epsilon}_0$, $\tilde{\epsilon}_0 \geq 0$, $f(x) \geq l_0(x) := \tilde{f}_0 + \langle y^0, x - x^0 \rangle$, $\forall x \in \mathfrak{R}^n$. Set $z^0 = x^0$, $f_0 := \tilde{f}_0$, and $k := 0$.

Step 1. Add the new piece to the approximation of f . Set

$$B_k^c := B_{k-1}^c \cup \{l_k(x)\}, \quad l_k(x) := \tilde{f}_k + \langle y^k, x - z^k \rangle,$$

$$I_k^c := \{0 \leq i \leq k \mid l_i(x) \in B_k^c\}.$$

Step 2. Minimize the regularized cutting plane approximation of f . Choose $\gamma_k > 0$ and compute z^{k+1} as the solution of

$$\min_x \{ \varphi_k(x) + (\gamma_k/2) \|x - x^k\|^2 \mid x \in \mathfrak{R}^n \}, \quad (7)$$

where

$$\varphi_k(x) := \max \{ \max \{ l_i(x) \mid i \in I_k^c \}, \max \{ l_k^a(x) \mid i \in I_k^a \} \}. \quad (8)$$

Step 3. Compute the value of the predicted descent,

$$\delta_k := \tilde{f}_k - \varphi_k(z^{k+1}) - (\gamma_k/2) \|z^{k+1} - x^k\|^2. \quad (9)$$

Step 4. Stopping test. Stop if

$$\delta_k + 2 \max \{ \tilde{\epsilon}_i \mid i \in I_k^c \} \leq \Delta_{\text{opt}}. \quad (10)$$

Otherwise, go to Step 5.

Step 5. Approximate the values of f and its subgradient at z^{k+1} . Compute some $\tilde{f}_{k+1} \in \mathfrak{R}$ and $y^{k+1} \in \mathfrak{R}^n$ such that

$$f(z^{k+1}) \geq \tilde{f}_{k+1} \geq f(z^{k+1}) - \tilde{\epsilon}_{k+1}, \quad \tilde{\epsilon}_{k+1} \geq 0, \quad (11a)$$

$$f(x) \geq l_{k+1}(x) := \tilde{f}_{k+1} + \langle y^{k+1}, x - z^{k+1} \rangle, \quad \forall x \in \mathfrak{R}^n. \quad (11b)$$

Step 6. Descent Test. If

$$\tilde{f}_k - \tilde{f}_{k+1} - \tilde{\epsilon}_{k+1} \geq \sigma \delta_k, \quad (12)$$

set $x^{k+1} := z^{k+1}$, $\tilde{f}_{k+1} := \tilde{f}_{k+1}$, $K_d := K_d \cup \{k+1\}$,

and go to Step 8.

Step 7. Null step. Set

$$x^{k+1} := x^k \quad \text{and} \quad \tilde{f}_{k+1} := \tilde{f}_k.$$

Step 8. Managing the complexity of the approximation of f . If $|B_k^c| + |B_k^a| < B_{\max}$, then set $B_{k+1}^a := B_k^a$ and go to Step 9. Otherwise, choose some $C_k \subset B_k^c \cup B_k^a$ such that $|C_k| \geq 2$, $l_{k_0}(x) \notin C_k$, where $k_0 = \max\{i | i \in K_d\}$. Set

$$B_k^c := B_k^c \setminus C_k, \quad B_{k+1}^a := (B_k^a \setminus C_k) \cup \{l_k^a(x)\},$$

$$l_k^a(x) := \varphi_k(z^{k+1}) + \gamma_k \|z^{k+1} - x^k\|^2 + \gamma_k \langle x^k - z^{k+1}, x - x^k \rangle,$$

$$I_{k+1}^a := \{0 \leq i \leq k+1 | l_i^a(x) \in B_{k+1}^a\}.$$

Step 9. Set $k := k+1$ and go to Step 1.

Note that the value of δ_k given by (9) is not necessarily nonnegative in the given setting. Another observation is that the stopping condition (10) is also somewhat unusual. According to Proposition 3.1, we could alternatively just use the more traditional (4), (5). However, this would require already at this stage all the objects defined in Lemma 3.1 below. Thus, we prefer the more compact form of (10).

3. Convergence Properties

We start with the dual characterization of the solution of (7), which is similar to what can be found in the literature on bundle methods; e.g. see Ref. 3 (Chapter 7). Here, we adapt it to our purposes and give a proof for completeness, as some of the expressions will be needed later on in any case.

Let $(\alpha_i, v^i, u^i) \in \mathfrak{R} \times \mathfrak{R}^n \times \mathfrak{R}^n$, $i = 1, \dots, |B_k^c| + |B_k^a|$ be the data defining φ_k in (8). We shall use the representation

$$\alpha_i + \langle v^i, x - u^i \rangle = \begin{cases} l_i(x) \in B_k^c, & i \in I_k^c, \\ l_i^a(x) \in B_k^a, & i \in I_k^a. \end{cases} \quad (13)$$

Note that, formally speaking, there are overlapping indices in (13), since $I_k^c \cap I_k^a \neq \emptyset$. While keeping in mind this detail, we prefer to use the

representation (13), as it does not lead to any real confusion, while simplifying notation. Let I_k be the union of all the indices defining φ_k [$|I_k| = |B_k^c| + |B_k^a|$], and consider the quadratic program

$$\max_{\lambda} \left\{ - (1/2\gamma_k) \left\| \sum_{i \in I_k} \lambda_i v^i \right\|^2 + \sum_{i \in I_k} \lambda_i (\alpha_i + \langle v^i, x^k - u^i \rangle) \mid \lambda \geq 0, \sum_{i \in I_k} \lambda_i = 1 \right\}, \quad (14)$$

where $\lambda \in \Re^{|I_k|}$.

Lemma 3.1. Let $\tilde{\lambda}^k$ be a solution of (14) and denote

$$d^k := \sum_{i \in I_k} \tilde{\lambda}_i^k v^i.$$

It holds that

$$z^{k+1} = x^k - (1/\gamma_k) d^k, \quad (15)$$

$$d^k \in \partial \varphi_k(z^{k+1}), \quad (16)$$

$$d^k \in \partial_{\epsilon_k} f(x^k), \quad (17)$$

where

$$\epsilon_k = \epsilon_k^c + \epsilon_k^a \geq 0,$$

with

$$\epsilon_k^c = \sum_{i \in I_k^c} \tilde{\lambda}_i^k [f(x^k) - f(z^i) - \langle y^i, x^k - z^i \rangle + \tilde{\epsilon}_i] \geq 0,$$

$$\epsilon_k^a = \sum_{i \in I_k^a} \tilde{\lambda}_i^k [f(x^k) - \varphi_i(z^{i+1}) - (1/\gamma_i) \|d^i\|^2 - \langle d^i, x^k - x^i \rangle] \geq 0.$$

Proof. Problem (7) is equivalent to the convex quadratic program

$$\min_{x, t} \{ t + (\gamma_k/2) \|x - x^k\|^2 \mid l_i(x) \leq t, i \in I_k^c, l_i^a(x) \leq t, i \in I_k^a \}.$$

The dual of this program (e.g. Ref. 21) is

$$\max_{x, t, \lambda} \left\{ t + (\gamma_k/2) \|x - x^k\|^2 + \sum_{i \in I_k} \lambda_i (\alpha_i + \langle v^i, x^k - u^i \rangle - t), \right. \\ \left. \gamma_k (x - x^k) + \sum_{i \in I_k} \lambda_i v^i = 0, \sum_{i \in I_k} \lambda_i = 1, \lambda \geq 0 \right\}, \quad (18)$$

with $(\alpha_i, v^i, u^i) \in \Re \times \Re^n \times \Re^n$, $i \in I_k$, defined according to (13). It can be seen easily that problem (18) reduces to (14) after eliminating the variables x and t via some simple transformations.

Assertion (15) of the lemma is now evident from the constraints in (18), using also the uniqueness of z^{k+1} and strong duality. Then, assertion (16) follows from the optimality condition for (7). We proceed to prove the last assertion.

We verify first by induction that

$$f(x) \geq \varphi_k(x), \quad \forall x \in \mathbb{R}^n, \forall k. \quad (19)$$

For $k = 0$, this is obvious, because $\varphi_0(x) = l_0(x)$. Suppose now that (19) holds for some index k . If $|B_k^c| + |B_k^a| < B_{\max}$, then

$$\varphi_{k+1}(x) = \max\{\varphi_k(x), l_{k+1}(x)\} \leq f(x),$$

by (11) and (19). If $|B_k^c| + |B_k^a| = B_{\max}$, then

$$\varphi_{k+1}(x) \leq \max\{\varphi_k(x), l_k^a(x), l_{k+1}(x)\} \leq f(x),$$

because

$$l_k^a(x) = \varphi_k(z^{k+1}) + \langle d^k, x - z^{k+1} \rangle \leq \varphi_k(x),$$

by (16). Hence, (19) holds.

By convex programming duality (e.g. Ref. 21), the optimal values of (7) and (14) are equal, from which we obtain that

$$\varphi_k(z^{k+1}) = -(1/\gamma_k) \|d^k\|^2 + \sum_{i \in I_k} \tilde{\lambda}_i^k [\alpha_i + \langle v^i, x^k - u^i \rangle]. \quad (20)$$

By (19), we obtain further that, for all $x \in \mathbb{R}^n$,

$$\begin{aligned} f(x) &\geq \varphi_k(x) \geq \varphi_k(z^{k+1}) + \langle d^k, x - z^{k+1} \rangle \\ &= \langle d^k, x - x^k \rangle + \sum_{i \in I_k} \tilde{\lambda}_i^k [\alpha_i + \langle v^i, x^k - u^i \rangle] \\ &= f(x^k) + \langle d^k, x - x^k \rangle - \sum_{i \in I_k} \tilde{\lambda}_i^k [f(x^k) - \alpha_i - \langle v^i, x^k - u^i \rangle] \\ &\geq f(x^k) + \langle d^k, x - x^k \rangle \\ &\quad - \sum_{i \in I_k^c} \tilde{\lambda}_i^k [f(x^k) - f(z^i) - \langle y^i, x^k - z^i \rangle + \tilde{\epsilon}_i] \\ &\quad - \sum_{i \in I_k^a} \tilde{\lambda}_i^k [f(x^k) - \varphi_i(z^{i+1}) - (1/\gamma_i) \|d^i\|^2 - \langle d^i, x^k - x^i \rangle] \\ &= f(x^k) + \langle d^k, x - x^k \rangle - \epsilon_k^c - \epsilon_k^a, \end{aligned}$$

where the second inequality is by (16), the first equality is by (20) and (15), and the third inequality is by (13) and (11). The quantity ϵ_k^c is nonnegative, because $y^i \in \partial_{\tilde{\epsilon}_i} f(z^i)$, by (11). To prove (17), it remains to estimate ϵ_k^a .

We have that

$$\begin{aligned}
 \epsilon_k^a &= \sum_{i \in I_k^a} \tilde{\lambda}_i^k [f(x^k) - \varphi_i(z^{i+1}) - (1/\gamma_i) \|d^i\|^2 - \langle d^i, x^k - x^i \rangle] \\
 &= \sum_{i \in I_k^a} \tilde{\lambda}_i^k [f(x^k) - \varphi_i(z^{i+1}) - \langle d^i, x^k - z^{i+1} \rangle] \\
 &\geq \sum_{i \in I_k^a} \tilde{\lambda}_i^k [f(x^k) - \varphi_i(x^k)] \\
 &\geq 0,
 \end{aligned}$$

where the second equality is by (15), the first inequality is by (16), and the last is by (19). \square

The next result shows that the stopping test of Algorithm 2.1 guarantees the desired accuracy in the solution of problem (1).

Proposition 3.1. Suppose that, at some iteration k , the stopping criterion (10) is satisfied. Then, x^k is an approximate solution of (1) in the sense of (4), (5), where d^k is defined in Lemma 3.1.

Proof. Let

$$k_0 = \max\{i \mid i \in K_d\}$$

be the index of the last descent step before (10) was satisfied. By (9) and taking into account that $\tilde{f}_k = \tilde{f}_{k_0}$, we have that

$$\begin{aligned}
 \delta_k &= \tilde{f}_k - \varphi_k(z^{k+1}) - (1/2\gamma_k) \|d^k\|^2 \\
 &= \tilde{f}_{k_0} + (1/2\gamma_k) \|d^k\|^2 - \sum_{i \in I_k} \tilde{\lambda}_i^k [\alpha_i + \langle v^i, x^k - u^i \rangle] \\
 &\geq (1/2\gamma_k) \|d^k\|^2 + \sum_{i \in I_k} \tilde{\lambda}_i^k [f(x^{k_0}) - \alpha_i - \langle v^i, x^k - u^i \rangle] - \tilde{\epsilon}_{k_0},
 \end{aligned}$$

where the second equality is by (20) and the inequality is by (11). Since $x^{k_0} = x^k$, we obtain further

$$\begin{aligned}
 \sum_{i \in I_k^c} \tilde{\lambda}_i^k [f(x^{k_0}) - \alpha_i - \langle v^i, x^k - u^i \rangle] &= \sum_{i \in I_k^c} \tilde{\lambda}_i^k [f(x^k) - \tilde{f}_i - \langle y^i, x^k - z^i \rangle] \\
 &\geq \sum_{i \in I_k^c} \tilde{\lambda}_i^k [f(x^k) - f(z^i) - \langle y^i, x^k - z^i \rangle] \\
 &= \epsilon_k^c - \sum_{i \in I_k^c} \tilde{\lambda}_i^k \tilde{\epsilon}_i \\
 &\geq \epsilon_k^c - \max\{\tilde{\epsilon}_i \mid i \in I_k^c\},
 \end{aligned}$$

where the first inequality is by (11). Taking also into account that

$$\sum_{i \in I_k^a} \tilde{\lambda}_i^k [f(x^k) - \alpha_i - \langle v^i, x^k - u^i \rangle] = \epsilon_k^a,$$

we conclude that

$$\begin{aligned} \delta_k &\geq (1/2\gamma_k) \|d^k\|^2 + \epsilon_k^c + \epsilon_k^a - \tilde{\epsilon}_{k_0} - \max\{\tilde{\epsilon}_i | i \in I_k^c\} \\ &\geq (1/2\gamma_k) \|d^k\|^2 + \epsilon_k - 2 \max\{\tilde{\epsilon}_i | i \in I_k^c\}, \end{aligned}$$

where the last inequality follows from the fact that $k_0 \in I_k^c$, by Step 8 of Algorithm 2.1. The last relation implies (5) whenever (10) holds. \square

Given Proposition 3.1, to validate our claims, it remains to prove that the stopping test (10) is necessarily activated if the approximation precision satisfies condition (6).

Theorem 3.1. Suppose that the optimal value of (1) is finite and that, in Algorithm 2.1, $\gamma_{k+1} \geq \gamma_k$ and $\tilde{\epsilon}_{k+1} \leq \tilde{\epsilon}_k$ for all k sufficiently large, $k \notin K_d$. Then, the following two events are mutually exclusive:

- (i) Algorithm 2.1 performs an infinite number of iterations.
- (ii) Condition (6) holds.

Proof. Suppose that $\{x^k\}$ is infinite, which means that (10) is never satisfied, i.e.,

$$\delta_k + 2 \max\{\tilde{\epsilon}_i | i \in I_k^c\} > \Delta_{\text{opt}}, \quad \forall k. \quad (21)$$

Consider first the case where there is an infinite number of descent steps. For all $k \in K_d$, (12) and (11) imply that

$$\begin{aligned} \sigma \delta_k &\leq \tilde{f}_k - (\tilde{f}_{k+1} + \tilde{\epsilon}_{k+1}) \\ &\leq f(x^k) - f(x^{k+1}). \end{aligned}$$

Hence,

$$\sigma \sum_{k \in K_d} \delta_k \leq f(x^0) - \min_x \{f(x) | x \in \mathcal{R}^n\} < +\infty.$$

In particular, it follows that

$$\liminf_{k \in K_d} \delta_k \leq 0.$$

By (21), we then have that

$$\limsup_{k \in K_d} (\max\{\tilde{\epsilon}_i | i \in I_k^c\}) \geq \Delta_{\text{opt}}/2,$$

which contradicts (6).

Suppose now that the number of descent steps is finite and $k_0 = \max\{i | i \in K_d\}$ is the index of the last descent iterate, so that $x^k = x^{k_0}$, $\tilde{f}_k = \tilde{f}_{k_0}$ for all $k \geq k_0$. For all k large enough, say $k \geq k_1$, we have that

$$\begin{aligned} & \tilde{f}_k - \delta_k + (\gamma_k/2) \|z^{k+2} - z^{k+1}\|^2 \\ &= \varphi_k(z^{k+1}) + (\gamma_k/2) [\|z^{k+1} - x^{k_0}\|^2 + \|z^{k+2} - z^{k+1}\|^2] \\ &= \varphi_k(z^{k+1}) + (\gamma_k/2) [2\langle x^{k_0} - z^{k+1}, z^{k+2} - z^{k+1} \rangle + \|z^{k+2} - x^{k_0}\|^2] \\ &\leq \varphi_k(z^{k+1}) + \langle d^k, z^{k+2} - z^{k+1} \rangle + (\gamma_{k+1}/2) \|z^{k+2} - x^{k_0}\|^2, \end{aligned} \quad (22)$$

where the first equality is by (9) and the inequality follows from (15) and $\gamma_k \leq \gamma_{k+1}$. Furthermore, if

$$|B_k^c| + |B_k^a| < B_{\max},$$

then

$$\begin{aligned} \varphi_k(z^{k+1}) + \langle d^k, z^{k+2} - z^{k+1} \rangle &\leq \varphi_k(z^{k+2}) \\ &\leq \max\{\varphi_k(z^{k+2}), l_{k+1}(z^{k+2})\} \\ &= \varphi_{k+1}(z^{k+2}), \end{aligned}$$

where the first inequality is by (16). If

$$|B_k^c| + |B_k^a| = B_{\max},$$

then

$$\varphi_k(z^{k+1}) + \langle d^k, z^{k+2} - z^{k+1} \rangle = l_k^a(z^{k+2}) \leq \varphi_{k+1}(z^{k+2}).$$

Combining these two cases, (22) implies that

$$\begin{aligned} & \tilde{f}_k - \delta_k + (\gamma_k/2) \|z^{k+2} - z^{k+1}\|^2 \\ &\leq \varphi_{k+1}(z^{k+2}) + (\gamma_{k+1}/2) \|z^{k+2} - x^{k_0}\|^2 \\ &= \tilde{f}_{k+1} - \delta_{k+1}, \end{aligned}$$

where the equality is by (9). Hence, for $k \geq k_1$,

$$\delta_k \geq \delta_{k+1} + (\gamma_k/2) \|z^{k+2} - z^{k+1}\|^2 \geq \delta_{k+1}. \quad (23)$$

In particular, the sequence $\{\delta_k\}$ is nonincreasing (for $k \geq k_1$) and it is bounded below by (21) and (6). Therefore, it converges,

$$\bar{\delta} = \lim_{k \rightarrow \infty} \delta_k. \quad (24)$$

We next show that the sequence $\{z^k\}$ is bounded. We have that

$$\begin{aligned} & \bar{f}_k - \delta_k + (\gamma_k/2) \|z^{k+1} - x^{k_0}\|^2 \\ &= \varphi_k(z^{k+1}) + \gamma_k \|z^{k+1} - x^{k_0}\|^2 \\ &= \varphi_k(z^{k+1}) + \langle d^k, x^{k_0} - z^{k+1} \rangle \\ &\leq \varphi_k(x^{k_0}), \end{aligned}$$

where the inequality follows from (16). Hence, for k large enough ($k \geq k_1$), it holds that

$$\begin{aligned} \|z^{k+1} - x^{k_0}\|^2 &\leq (2/\gamma_k) [\delta_k + \varphi_k(x^{k_0}) - \bar{f}_{k_0}] \\ &\leq (2/\gamma_{k_1}) [\delta_{k_1} + f(x^{k_0}) - \bar{f}_{k_0}] \\ &\leq (2/\gamma_{k_1}) (\delta_{k_1} + \tilde{\epsilon}_{k_0}), \end{aligned}$$

where the second inequality is by (19) and by the monotonicity of $\{\gamma_k\}$ and $\{\delta_k\}$ (for $k \geq k_1$) and the last inequality is by (11). Thus, $\{z^k\}$ is bounded.

Since for $k \geq k_0$ the descent test (12) is never satisfied, we have that

$$\bar{f}_k - \tilde{f}_{k+1} - \tilde{\epsilon}_{k+1} < \sigma \delta_k.$$

Hence,

$$\begin{aligned} (1 - \sigma) \delta_k &< \tilde{f}_{k+1} + \tilde{\epsilon}_{k+1} - \bar{f}_k + \delta_k \\ &\leq f(z^{k+1}) + \tilde{\epsilon}_{k+1} - \varphi_k(z^{k+1}) - (\gamma_k/2) \|z^{k+1} - x^{k_0}\|^2 \\ &\leq f(z^{k+1}) - \tilde{f}_k + \varphi_k(z^k) - \varphi_k(z^{k+1}) + \tilde{\epsilon}_{k+1} \\ &\leq f(z^{k+1}) - f(z^k) + \tilde{\epsilon}_k + \varphi_k(z^k) - \varphi_k(z^{k+1}) + \tilde{\epsilon}_{k+1} \\ &\leq 2L \|z^{k+1} - z^k\| + 2 \max\{\tilde{\epsilon}_i | i \in I_k^c\}, \end{aligned}$$

where the second inequality is by (11) and (9), the third is by $\varphi_k(z^k) \geq l_k(z^k) = \tilde{f}_k$, the fourth by (11), and the last is by the Lipschitz continuity of f and φ_k , taking into account that $\{z^k\}$ is bounded.

Combining the latter relation with (23), we obtain that

$$\delta_k - \delta_{k+1} \geq (\gamma_k/8L^2) [(1 - \sigma) \delta_{k+1} - 2 \max\{\tilde{\epsilon}_i | i \in I_{k+1}^c\}]^2.$$

Since, by (24), $\delta_k - \delta_{k+1} \rightarrow 0$, and taking also into account that $\gamma_k \geq \gamma_{k_1} > 0$, the relation above implies that

$$\lim_{k \rightarrow \infty} \delta_k = \bar{\delta} = (2/(1 - \sigma)) \lim_{k \rightarrow \infty} (\max\{\tilde{\epsilon}_i | i \in I_k^c\}).$$

Passing onto the limit in (21), we then obtain

$$(2 + 2/(1 - \sigma)) \lim_{k \rightarrow \infty} (\max\{\tilde{\epsilon}_i | i \in I_k^c\}) \geq \Delta_{\text{opt}},$$

in contradiction with (6). □

As a final remark, we note that the presented analysis provides also a stability result for the proximal bundle algorithm. Indeed, instead of stating the question about how small should be the approximation errors $\tilde{\varepsilon}$ in order to guarantee that the iterates would eventually satisfy the given optimality tolerance Δ_{opt} , we could pose the following closely related questions. Given an upper bound $\varepsilon > 0$ for the approximation errors $\tilde{\varepsilon}$, i.e.,

$$\varepsilon \geq \limsup_k (\max\{\tilde{\varepsilon}_i | i \in I_k^c\}),$$

what are the convergence properties of the proximal bundle method? What kind of an approximate solution can be obtained? How does it depend on the value of ε ? Our results show that, using the bundle method, we are guaranteed to find an approximate solution of (1) in the sense of (4), (5), with

$$\Delta_{\text{opt}} = 2(2 - \sigma)\varepsilon/(1 - \sigma) + t, \quad t > 0 \text{ arbitrarily small.}$$

Our analysis applies even when the value of ε can be relatively large. One example where the latter is possible are problems with noisy data, where the bound on the magnitude of perturbations is known, but usually not controllable (in particular, it cannot be driven to zero).

References

1. KIWIEL, K. C., *Methods of Descent for Nondifferentiable Optimization*. Lecture Notes in Mathematics, Springer Verlag, Berlin, Germany, Vol. 1133, 1985.
2. HIRIART-URRUTY, J. B., and LEMARÉCHAL, C., *Convex Analysis and Minimization Algorithms*, Springer Verlag, Berlin, Germany, 1993.
3. BONNANS, J. F., GILBERT, J. C., LEMARÉCHAL, C., and SAGASTIZÁBAL, C. A., *Optimisation Numérique: Aspects Théoriques et Pratiques*, Springer Verlag, Berlin, Germany, 1997.
4. BERTSEKAS, D. P., and TSITSIKLIS, J. N., *Parallel and Distributed Computation*, Prentice-Hall, Englewood Cliffs, New Jersey, 1989.
5. LEMARÉCHAL, C., *Lagrangian Decomposition and Nonsmooth Optimization: Bundle Algorithm, Prox Iteration, Augmented Lagrangian*, in *Nonsmooth Optimization: Methods and Applications*, Edited by F. Giannessi, Gordon and Breach, Philadelphia, Pennsylvania, pp. 201–216, 1992.
6. BERTSEKAS, D. P., *Nonlinear Programming*, Athena Scientific, Belmont, Massachusetts, 1995.
7. SOLODOV, M. V., *Convergence Analysis of Perturbed Feasible Descent Methods*, Journal of Optimization Theory and Applications, Vol. 93, pp. 337–353, 1997.
8. SOLODOV, M. V., and ZAVRIEV, S. K., *Error Stability Properties of Generalized Gradient-Type Algorithms*, Journal of Optimization Theory and Applications, Vol. 98, pp. 663–680, 1998.

9. KIWIEL, K. C., *Approximations in Proximal Bundle Methods and Decomposition of Convex Programs*, Journal of Optimization Theory and Applications, Vol. 84, pp. 529–548, 1995.
10. HINTERMÜLLER, M., *A Proximal Bundle Method Based on Approximate Subgradients*, Computational Optimization and Applications, Vol. 20, pp. 245–266, 2001.
11. MILLER, S. A., *An Inexact Bundle Method for Solving Large Structured Linear Matrix Inequalities*, PhD Thesis, University of California, Santa Barbara, California, 2001.
12. KIWIEL, K. C., *A Method for Solving Certain Quadratic Programming Problems Arising in Nonsmooth Optimization*, IMA Journal of Numerical Analysis, Vol. 6, pp. 137–152, 1986.
13. KIWIEL, K. C., *Proximity Control in Bundle Methods for Convex Nondifferentiable Minimization*, Mathematical Programming, Vol. 46, pp. 105–122, 1990.
14. SCHRAMM, H., and ZOWE, J., *A Version of the Bundle Idea for Minimizing a Nonsmooth Function: Conceptual Idea, Convergence Analysis, Numerical Results*, SIAM Journal on Optimization, Vol. 2, pp. 121–152, 1992.
15. BONNANS, J. F., GILBERT, J. C., LEMARÉCHAL, C., and SAGASTIZÁBAL, C., *A Family of Variable-Metric Proximal-Point Methods*, Mathematical Programming, Vol. 68, pp. 15–47, 1995.
16. LEMARÉCHAL, C., and SAGASTIZÁBAL, C., *An Approach to Variable-Metric Bundle Methods*, System Modelling and Optimization, Lecture Notes in Control and Information Sciences, Edited by J. Henry and J. P. Yvon, Springer, Berlin, Germany, Vol. 197, pp. 144–162, 1994.
17. MIFFLIN, R., *A Quasi-Second-Order Proximal Bundle Algorithm*, Mathematical Programming, Vol. 73, pp. 51–72, 1996.
18. LEMARÉCHAL, C., and SAGASTIZÁBAL, C., *Variable-Metric Bundle Methods: From Conceptual to Implementable Forms*, Mathematical Programming, Vol. 76, pp. 393–410, 1997.
19. CHEN, X., and FUKUSHIMA, M., *Proximal Quasi-Newton Methods for Nondifferentiable Convex Optimization*, Mathematical Programming, Vol. 85, pp. 313–334, 1999.
20. LUKŠAN, L., and VLČEK, J., *Globally Convergent Variable-Metric Method for Convex Nonsmooth Unconstrained Optimization*, Journal of Optimization Theory and Applications, Vol. 102, pp. 593–613, 1999.
21. MANGASARIAN, O. L., *Nonlinear Programming*, McGraw-Hill, New York, New York, 1969.