

A New Sequential Optimality Condition for Constrained Nonsmooth Optimization

Elias Salomão Helou* Sandra A. Santos† Lucas E. A. Simões†

November 23, 2018

Abstract

We introduce a sequential optimality condition for locally Lipschitz constrained nonsmooth optimization, verifiable just using derivative information, and which holds even in the absence of any constraint qualification. The proposed sequential optimality condition is not only novel for nonsmooth problems, but brings new insights for the smooth case as well. We present a practical algorithm that generates iterates fulfilling the new necessary optimality condition. A main feature of the devised algorithm is to allow a stronger control over the infeasibility of the iterates than usually obtained by exact penalty strategies, ensuring theoretical and practical advantages. Illustrative numerical experiments highlight the potentialities of the algorithm.

Keywords nonsmooth nonconvex optimization; constrained optimization; sequential optimality condition; constraint qualification

AMS Subject Classifications 90C26, 90C30, 90C46

1 Introduction

We consider constrained nonsmooth optimization problems of the form

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \\ \text{s.t. } \mathbf{c}(\mathbf{x}) \leq \mathbf{0}, \end{aligned} \tag{P}$$

where both $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $\mathbf{c} : \mathbb{R}^n \rightarrow \mathbb{R}^p$ are locally Lipschitz continuous functions. Equality constraints of the type $\mathbf{h}(\mathbf{x}) = 0$ can also be easily incorporated in our framework (see Remark 1 for more details).

A common way of solving constrained optimization problems is to turn (P) into an unconstrained minimization by penalizing points $\mathbf{x} \in \mathbb{R}^n$ that violate the constraints, i.e., one seeks to find a solution to

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) + \rho z(\mathbf{x}), \tag{Unc-P}$$

where ρ is a positive real number and $z : \mathbb{R}^n \rightarrow \mathbb{R}$ is a continuous function satisfying

$$\begin{cases} z(\mathbf{x}) > 0, & \text{if } \mathbf{c}(\mathbf{x}) > \mathbf{0} \\ z(\mathbf{x}) = 0, & \text{if } \mathbf{c}(\mathbf{x}) \leq \mathbf{0}. \end{cases}$$

*Institute of Mathematical Sciences and Computation - University of São Paulo. São Carlos - SP, Brazil. (elias@icmc.usp.br)

†Department of Applied Mathematics - University of Campinas. Campinas - SP, Brazil. (sandra@ime.unicamp.br, simoes.lea@gmail.com)

Under certain conditions, and considering ρ to be large enough, the unconstrained minimization (Unc-P) and the original optimization (P) are equivalent [12]. This approach is known as *exact penalization* [28, 50] and it can be used to derive most of the theory behind constrained optimization for finite dimensions [13].

For the majority of the optimization problems, the equivalence between (P) and (Unc-P) can only be achieved for finite ρ when z is a nondifferentiable function. However, when both objective and constraint functions are smooth in the entire domain, one usually does not want to trade the constrained smooth problem by a nonsmooth unconstrained one. Hence, many alternatives have been proposed over the years in order to overcome the nondifferentiable nature of (Unc-P). For instance, the Augmented Lagrangian [1, 20] and the Sequential Quadratic Programming (SQP) [9, 46] are widely known methods. However, when (P) already presents nonsmoothness, the nondifferentiability of z usually does not introduce any additional difficulty. For this reason, many methods developed for solving constrained nonsmooth optimization problems are based on exact penalty functions [22, 23, 38, 39, 48].

There is a strong connection between constraint qualifications (CQ) [5, 12, 46] for problem (P) and the existence of a finite penalty parameter that makes (Unc-P) equivalent to (P). Moreover, convergence results for exact penalization methods (and, more generally, for the majority of the optimization algorithms) are based on different kinds of CQs. However, checking the validity of constraint qualifications is not an easy task, which may justify the fact that most of the algorithms are not designed to test any kind of CQ even when theoretical convergence of the method relies on such conditions. As a result, practical necessary optimality conditions for (P) that do not depend on constraint qualifications are of great importance for establishing theoretically sound stopping criteria for any optimization method.

The *approximate Karush-Kuhn-Tucker* (AKKT) and the *complementary approximate Karush-Kuhn-Tucker* (CAKKT) conditions [2, 4] present themselves as reliable necessary optimality conditions for smooth optimization problems. Looking at auxiliary functions that approximate the exact penalty approach, the authors of both studies show that any solution \mathbf{x}^* of the smooth optimization problem must satisfy the following sequential optimality condition

$$\lim_{k \rightarrow \infty} \left\| \nabla f(\mathbf{x}^k) + \sum_{i=1}^p \nabla c_i(\mathbf{x}^k) \mu_i^k \right\| = 0, \quad (1)$$

where $\|\mathbf{x}\|$ is the Euclidean norm of $\mathbf{x} \in \mathbb{R}^n$, $\{\mathbf{x}^k\} \subset \mathbb{R}^n$ is a sequence converging to \mathbf{x}^* , and $\{\boldsymbol{\mu}^k\} \subset \mathbb{R}_+^p$ is a sequence of vectors whose components μ_i^k must satisfy $\lim_{k \rightarrow \infty} \min\{-c_i(\mathbf{x}^k), \mu_i^k\} = 0$, where $c_i : \mathbb{R}^n \rightarrow \mathbb{R}$ are the components of $\mathbf{c} : \mathbb{R}^n \rightarrow \mathbb{R}^p$ for each $i \in \{1, \dots, p\}$. A natural attempt to generalize (1) to the nonsmooth case is to use the subdifferential concept, i.e., to consider a sequence of elements

$$\{\mathbf{v}^k\} \subset \left(\partial f(\mathbf{x}^k) + \sum_{i=1}^p \partial c_i(\mathbf{x}^k) \mu_i^k \right), \quad \text{with } \lim_{k \rightarrow \infty} \|\mathbf{v}^k\| = 0, \quad (2)$$

where $\partial g(\mathbf{x})$ stands for the Clarke subdifferential set of the function g at the point \mathbf{x} [18]. Unfortunately, appropriate vectors from the subdifferential satisfying these conditions are, in general, not expected to be computed by numerical algorithms, which invalidates the use of (2) as a practical stopping criteria.

Using the fact that many real nonsmooth optimization problems are described by functions that are continuously differentiable in a full-measure subset of \mathbb{R}^n , the authors of [27] present a necessary optimality condition that relies upon the fact that, in many cases, the sets $\partial f(\mathbf{x}^k)$ and $\partial c_i(\mathbf{x}^k)$ can be traded by $\nabla f(\hat{\mathbf{x}}^k)$ and $\nabla c_i(\hat{\mathbf{x}}^k)$, where $\hat{\mathbf{x}}^k$ is a point sufficiently close to \mathbf{x}^k . However, if no constraint qualification is assumed, one must have $\hat{\mathbf{x}}^k = \mathbf{x}^k$, which brings back the same issues discussed in the previous paragraph.

In this paper, we propose a new sequential optimality condition for constrained nonsmooth optimization problems that allows the user to work with derivatives even in the absence of any CQ. To

show that our necessary optimality condition is practical, we present an algorithm that is proven to generate a sequence of points that satisfies such conditions. A main feature of the proposed algorithm is that it presents a stronger control over the infeasibility of the iterates than usual exact penalty methods, which has important practical consequences.

One of the main challenges of exact penalty methods is the selection of the initial value of the penalty parameter. In case this value is chosen to be too large, the method may privilege the feasible region in detriment of optimality, which, in turn, can cause very short steps toward the optimal solution and/or may cause the penalized problem to be ill-conditioned. On the other hand, if the penalty parameter value is much smaller than the magnitude of the objective function at infeasible points, the method may rapidly be attracted by unconstrained minimizers that possess very low function values, preventing the user to obtain a successful solution for reasonable values of ρ . Such a phenomenon is called *greediness* [7, 17], and may occur even if the user starts the method at a feasible point of the optimization problem.

Unlike some existing penalty methods, for which there is no simple way to compute a suitable initial value for the penalty parameter to easily confine the iterates into an almost feasible region, our algorithm allows the user to set a tolerance target value $\xi > 0$ for infeasibility. Outside this tolerance region, the method is indifferent to the objective function, which ensures the control of the infeasibility even at the initial iterations. As a consequence, it prevents our method to suffer from the greediness phenomenon.

The outline of the paper is as follows. Section 2 presents theoretical results supporting that every local minimizer of (P) must fulfill our proposed sequential optimality conditions. Section 3 shows the relation between our sequential optimality conditions and the AKKT and CAKKT conditions in the case where the objective and constraint functions of (P) are all smooth. In Section 4, we state a general algorithm that produces a sequence of iterates fulfilling the proposed sequential optimality conditions. Section 5 brings numerical results that illustrate some of the important properties of our method. Finally, we leave Section 6 for the conclusions of this study.

During the entire manuscript, the following notations will be frequently used:

- $\|\cdot\|$ is the Euclidean norm;
- $\mathcal{B}(\mathbf{x}, \epsilon)$ is the Euclidean closed ball with center at \mathbf{x} and radius ϵ ;
- $\mathcal{P}(\mathbf{v} | \mathcal{X})$ is the Euclidean projection of \mathbf{v} onto \mathcal{X} ;
- \mathbb{R}_+ is the set of all nonnegative real numbers;
- \mathbb{R}_+^* is the set of all strictly positive real numbers;
- $\mathbf{v}_+ := \max\{v, 0\}$. If v is a vector, then \mathbf{v}_+ is also a vector where every entry is taken as the maximum between the respective entry of \mathbf{v} and zero;
- $\text{conv } \mathcal{X}$ denotes the convex hull of \mathcal{X} ;
- $\text{cl } \mathcal{X}$ denotes the closure of \mathcal{X} ;
- $A + B$ represents the set $\{x + y : x \in A, y \in B\}$;
- $A \cdot B$ represents the set $\{xy : x \in A, y \in B\}$.

2 Establishing new necessary optimality conditions

We start this section by introducing the concept of (Goldstein) ϵ -subdifferential set for any locally Lipschitz continuous function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ [31].

Definition 2.1 (ϵ -Subdifferential set, ϵ -subgradient, ϵ -stationary point). *The ϵ -subdifferential set of f at x is given by*

$$\partial_\epsilon f(\mathbf{x}) := \text{conv } \partial f(\mathcal{B}(\mathbf{x}, \epsilon)),$$

where $\text{conv } \mathcal{S}$ is the convex hull of the set \mathcal{S} and $\mathcal{B}(\mathbf{x}, \epsilon) := \{\mathbf{y} : \|\mathbf{y} - \mathbf{x}\| \leq \epsilon\}$. Any $\mathbf{v} \in \partial_\epsilon f(\mathbf{x})$ is known as an ϵ -subgradient of f at x . Moreover, if $\mathbf{0} \in \partial_\epsilon f(\mathbf{x})$, then we say that x is an ϵ -stationary point for f .

One of the greatest advantages of looking at ϵ -subdifferential sets instead of subdifferential sets is the possibility of seeing $\partial_\epsilon f(\mathbf{x})$ as a convex hull of the derivatives of f nearby x . Indeed, for any locally Lipschitz continuous function f and any full-measure subset \mathcal{D}^f of \mathbb{R}^n such that f is differentiable at any point in \mathcal{D}^f – this subset always exists due to Rademachers theorem [30, Theorem 3.1.6] – the following relations hold (see [14, 40])

$$\mathcal{G}_\epsilon^f(\mathbf{x}) \subset \partial_\epsilon f(\mathbf{x}) \quad \text{and} \quad \partial_{\epsilon_1} f(\mathbf{x}) \subset \mathcal{G}_{\epsilon_2}^f(\mathbf{x}) \quad (0 \leq \epsilon_1 < \epsilon_2), \quad (3)$$

where $\mathcal{G}_\epsilon^f(\mathbf{x}) := \text{cl conv } \nabla f(\mathcal{B}(\mathbf{x}, \epsilon) \cap \mathcal{D}^f)$, with $\text{cl } \mathcal{S}$ being the closure of the set \mathcal{S} . Recent advances on practical tools [14, 15, 16] to approximate $\mathcal{G}_\epsilon^f(\mathbf{x})$ allow us to consider a new sequential optimality condition of practical use. The subgradients involved in this new necessary optimality condition are those associated with the objective function and the constraints related to the following index set

$$\mathcal{I}_\epsilon(\mathbf{x}) := \{i : \exists \mathbf{y} \in \mathcal{B}(\mathbf{x}, \epsilon) \text{ with } c_i(\mathbf{y}) \geq 0\}.$$

Definition 2.2 (weak ϵ -Approximate Nonsmooth Optimality Condition). *A feasible point $\mathbf{x}^* \in \mathbb{R}^n$ of (P) is said to satisfy the weak ϵ -Approximate Nonsmooth Optimality Condition (weak ϵ -ANOC) if there exist sequences $\{\mathbf{x}^k\} \subset \mathbb{R}^n$, $\{\epsilon_k\} \subset \mathbb{R}_+^*$, where \mathbb{R}_+^* is the set of positive real numbers, $\{\mathbf{v}^k\} \subset \mathbb{R}^n$ and $\{\mu^k\} \subset \mathbb{R}_+^p$ such that $\mathbf{x}^k \rightarrow \mathbf{x}^*$, $\epsilon_k \downarrow 0$ and $\|\mathbf{v}^k\| \rightarrow 0$, where*

$$\mathbf{v}^k \in \left(\mathcal{G}_{\epsilon_k}^f(\mathbf{x}^k) + \sum_{i=1}^p \mu_i^k \mathcal{G}_{\epsilon_k}^{c_i}(\mathbf{x}^k) \right) \quad \text{and} \quad i \notin \mathcal{I}_{\epsilon_k}(\mathbf{x}^k) \Rightarrow \mu_i^k = 0. \quad (4)$$

Notice that the above definition does not impose any control over the speed of the convergence $\epsilon_k \downarrow 0$. Not establishing a relation between the sequences $\{\epsilon_k\}$ and $\{\mathbf{x}^k\}$ allows $\epsilon_k \gg \|\mathbf{x}^k - \mathbf{x}^*\|$, which, in turn, means that the weak ϵ -ANOC may depart too much from the necessary condition that at least one generalized derivative of f at \mathbf{x}^k must be close to a linear combination of the subgradients of the active constraints at \mathbf{x}^k . Fig. 1 exemplifies this issue, showing that subgradients of far away inactive constraints may be considered in the linear combination.

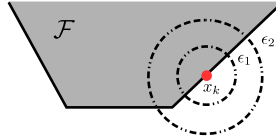


Figure 1: Illustration of how the value of ϵ in the ϵ -subdifferential can influence in the activeness of the constraints.

Using the exact penalization approach for problem (P) and a similar reasoning employed by the authors in [4, Theorem 3.3], one can infer the weak ϵ -ANOC as a necessary optimality condition for nonsmooth problems. However, this strategy does not clarify on the question of how fast the sequence $\{\epsilon_k\}$ must go to zero. To overcome this issue, we have applied a new penalization strategy that has its roots in [8, Section 4.1].

The inspiring idea behind this new penalization is to trade the original nonsmooth problem by the minimization of a discontinuous function $\Theta_\rho : \mathbb{R}^n \rightarrow \mathbb{R}$,

$$\begin{cases} \Theta_\rho(\mathbf{x}) \geq 0 & \text{if } \mathbf{c}(\mathbf{x}) > \mathbf{0} \\ \Theta_\rho(\mathbf{x}) = f(\mathbf{x}) - \rho & \text{otherwise,} \end{cases}$$

where $\rho \in \mathbb{R}$ plays a different role than the one presented in the exact penalization function. Here, the parameter ρ is any constant ensuring that $f(\mathbf{x}^*) - \rho < 0$, with \mathbf{x}^* being a solution of (P). Like the strategy of solving constrained smooth optimization problems by the use of an exact penalty function, where a degree of difficulty is added by introducing a nondifferentiable term in order to avoid a constrained minimization, this approach also inserts one extra difficulty to the nonsmoothness of the original problem by using a discontinuous function. Therefore, we avoid using this idea directly.

Theorem 2.1 ahead shows that the weak ϵ -ANOC is, indeed, a necessary optimality condition for problem (P) and also elucidates the matter of the speed of the convergence $\epsilon_k \downarrow 0$. Its proof relies on a sequence $\{\Psi_{\xi_k, \rho}\}$ of nonsmooth continuous functions $\Psi_{\xi_k, \rho} : \mathbb{R}^n \rightarrow \mathbb{R}$ approximating the discontinuous function Θ_ρ . Given a scalar $\xi \in \mathbb{R}_+^*$, we define

$$\Psi_{\xi, \rho}(\mathbf{x}) := \max \left\{ 1 - \frac{\|\mathbf{c}(\mathbf{x})_+\|_1}{\xi}, 0 \right\} [f(\mathbf{x}) - \rho] + \|\mathbf{c}(\mathbf{x})_+\|_1. \quad (5)$$

It is worth noticing that, since f and \mathbf{c} are locally Lipschitz continuous functions, the map $\Psi_{\xi, \rho}$ is also locally Lipschitz continuous.

Preceding the result that establishes the weak ϵ -ANOC as a necessary optimality condition, we present two lemmas that will facilitate the proof of Theorem 2.1.

Lemma 2.1. *Let $\mathbf{x} \in \mathbb{R}^n$, $(\epsilon, \xi) \in \mathbb{R}_+^* \times \mathbb{R}_+^*$, $\rho \in \mathbb{R}$ with $\rho > f(\mathbf{x})$, and $\mathbf{v} \in \partial_\epsilon \Psi_{\xi, \rho}(\mathbf{x})$. Then, if $\epsilon > 0$ is such that $\mathbf{y} \in \mathcal{B}(\mathbf{x}, \epsilon) \Rightarrow \rho \geq f(\mathbf{y})$, there exists $\{\mathbf{x}^j\}_{j=1}^{n+1} \subset \mathcal{B}(\mathbf{x}, \epsilon)$ and $\boldsymbol{\lambda} \in \mathbb{R}_+^{n+1}$, with $\mathbf{e}^T \boldsymbol{\lambda} = 1$, where $\mathbf{e} := (1, 1, \dots, 1)^T$, such that*

$$\mathbf{v} \in \sum_{j=1}^{n+1} \lambda_j \left(\max \left\{ 1 - \frac{\|\mathbf{c}(\mathbf{x}^j)_+\|_1}{\xi}, 0 \right\} \partial f(\mathbf{x}^j) + \sigma_j \partial \|\mathbf{c}(\mathbf{x}^j)_+\|_1 \right), \quad (6)$$

where $\sigma_j \geq 1$, $j \in \{1, \dots, n+1\}$.

Proof. Recalling that $\partial_\epsilon \Psi_{\xi, \rho}(\mathbf{x})$ is a convex set, it follows from Carathéodory's Theorem [49, Theorem 2.29] that there exists $\{\mathbf{s}^j\}_{j=1}^{n+1} \subset \partial \Psi_{\xi, \rho}(\mathcal{B}(\mathbf{x}, \epsilon))$ and $\boldsymbol{\lambda} \in \mathbb{R}_+^{n+1}$, with $\mathbf{e}^T \boldsymbol{\lambda} = 1$, such that $\mathbf{v} = \sum_{j=1}^{n+1} \lambda_j \mathbf{s}^j$. Because of [18, Theorem 2.3.13], [18, Theorem 2.3.9], and [18, Proposition 2.3.3], we know that $\partial \Psi_{\xi, \rho}(\mathbf{x})$ is a subset of

$$\max \left\{ 1 - \frac{\|\mathbf{c}(\mathbf{x})_+\|_1}{\xi}, 0 \right\} \partial f(\mathbf{x}) + \left(1 + [\rho - f(\mathbf{x})] \text{conv} \left\{ \frac{1}{\xi}, 0 \right\} \right) \cdot \partial \|\mathbf{c}(\mathbf{x})_+\|_1.$$

So, since, by hypothesis, $\rho - f(\mathbf{x}^j) > 0$, $j \in \{1, \dots, n+1\}$, it follows that there exists $\{\mathbf{x}^j\}_{j=1}^{n+1} \subset \mathcal{B}(\mathbf{x}, \epsilon)$ such that (6) holds with $1 \leq \sigma_j \in \left(1 + [\rho - f(\mathbf{x}^j)] \text{conv} \left\{ \frac{1}{\xi}, 0 \right\} \right)$. \square

The next result presents sufficient conditions for the feasible point \mathbf{x}^* to satisfy the weak ϵ -ANOC.

Lemma 2.2. *Suppose \mathbf{x}^* is a feasible point of problem (P), $\{\mathbf{x}^k\}$ is a sequence converging to \mathbf{x}^* , $\{\epsilon_k\}$ and $\{\xi_k\}$ are both real-valued sequences with $\epsilon_k \downarrow 0$, $\xi_k \downarrow 0$ and $\epsilon_k/\xi_k \rightarrow 0$, and ρ is a real value satisfying $f(\mathbf{x}^*) - \rho < 0$. If $\lim_{k \rightarrow \infty} \|\mathbf{c}(\mathbf{x}^k)_+\|_1/\xi_k$ exists and it is strictly less than one, and there exists a sequence $\{\mathbf{r}^k\} \subset \mathbb{R}^n$ such that $\|\mathbf{r}^k\| \rightarrow 0$ and $\mathbf{r}^k \in \partial_{\epsilon_k} \Psi_{\xi_k, \rho}(\mathbf{x}^k)$ for all $k \in \mathbb{N}$, then \mathbf{x}^* satisfies the weak ϵ -ANOC.*

Proof. Let us choose a sequence $\{\zeta_k\}$ satisfying $\zeta_k < \epsilon_k$ for all $k \in \mathbb{N}$. Recalling Lemma 2.1, it follows that, for all large enough k , there exists $\{\mathbf{x}^{k,j}\}_{j=1}^{n+1} \subset \mathcal{B}(\mathbf{x}^k, \zeta_k)$ such that

$$\mathbf{r}^k \in \sum_{j=1}^{n+1} \lambda_j^k \left(\max \left\{ 1 - \frac{\|\mathbf{c}(\mathbf{x}^{k,j})_+\|_1}{\xi_k}, 0 \right\} \partial f(\mathbf{x}^{k,j}) + \sigma_j^k \partial \|\mathbf{c}(\mathbf{x}^{k,j})_+\|_1 \right),$$

with $\sigma_j^k \geq 1$, $j \in \{1, \dots, n+1\}$. Therefore, for each j , there must exist $\mathbf{s}_f^{k,j} \in \partial f(\mathbf{x}^{k,j})$ and $\mathbf{s}_{c_i}^{k,j} \in \partial c_i(\mathbf{x}^{k,j})$ for each i , such that

$$\mathbf{r}^k = \sum_{j=1}^{n+1} \lambda_j^k \left(\max \left\{ 1 - \frac{\|\mathbf{c}(\mathbf{x}^{k,j})_+\|_1}{\xi_k}, 0 \right\} \mathbf{s}_f^{k,j} + \sigma_j^k \sum_{i=1}^p \Delta_i(\mathbf{x}^{k,j}) \mathbf{s}_{c_i}^{k,j} \right), \quad (7)$$

where

$$\Delta_i(\mathbf{x}) \in \begin{cases} \{1\}, & \text{if } c_i(\mathbf{x}) > 0 \\ \text{conv}\{1, 0\}, & \text{if } c_i(\mathbf{x}) = 0 \\ \{0\}, & \text{if } c_i(\mathbf{x}) < 0 \end{cases}.$$

Then, defining

$$l_k^f = \sum_{j=1}^{n+1} \lambda_j^k \max \left\{ 1 - \frac{\|\mathbf{c}(\mathbf{x}^{k,j})_+\|_1}{\xi_k}, 0 \right\} \quad \text{and} \quad l_k^{c_i} = \sum_{j=1}^{n+1} \lambda_j^k \sigma_j^k \Delta_i(\mathbf{x}^{k,j}), \quad (8)$$

one can see, for large values of k , that l_k^f is strictly positive due to the fact that $\lim_{k \rightarrow \infty} \|\mathbf{c}(\mathbf{x}^k)_+\|_1 / \xi_k$ is strictly less than one, $\|\mathbf{x}^k - \mathbf{x}^{k,j}\| \leq \epsilon_k$, and $\epsilon_k / \xi_k \rightarrow 0$. Hence, for k large enough, we have, for all i and j , that

$$0 \leq \tau_{k,j}^f := \frac{\lambda_j^k \max \left\{ 1 - \frac{\|\mathbf{c}(\mathbf{x}^{k,j})_+\|_1}{\xi_k}, 0 \right\}}{l_k^f} \quad \text{and} \quad 0 \leq \tau_{k,j}^{c_i} := \begin{cases} \frac{\lambda_j^k \sigma_j^k \Delta_i(\mathbf{x}^{k,j})}{l_k^{c_i}}, & \text{if } l_k^{c_i} > 0 \\ 0, & \text{if } l_k^{c_i} = 0 \end{cases}.$$

Additionally, the following holds: $\sum_{j=1}^{n+1} \tau_{k,j}^f = 1$ and

$$\sum_{j=1}^{n+1} \tau_{k,j}^{c_i} = \begin{cases} 1, & \text{if } l_k^{c_i} > 0 \\ 0, & \text{if } l_k^{c_i} = 0 \end{cases} \quad \text{for all } i \in \{1, \dots, p\}.$$

This together with (7) imply that $\mathbf{r}^k \in \left(l_k^f \partial_{\zeta_k} f(\mathbf{x}^k) + \sum_{i=1}^p l_k^{c_i} \partial_{\zeta_k} c_i(\mathbf{x}^k) \right)$. So, choosing $\mathbf{v}^k := \frac{1}{l_k^f} \mathbf{r}^k$, we obtain $\mathbf{v}^k \in \left(\partial_{\zeta_k} f(\mathbf{x}^k) + \sum_{i=1}^p \frac{l_k^{c_i}}{l_k^f} \partial_{\zeta_k} c_i(\mathbf{x}^k) \right)$. Now, notice that, because of the way $l_k^{c_i}$ and Δ_i were defined, it must follow $i \notin \mathcal{I}_{\zeta_k}(\mathbf{x}^k) \Rightarrow \mu_i^k := \frac{l_k^{c_i}}{l_k^f} = 0$. Moreover, since $\|\mathbf{r}^k\| \rightarrow 0$, it yields $\|\mathbf{v}^k\| \rightarrow 0$ as well. So,

$$\mathbf{v}^k \in \left(\partial_{\zeta_k} f(\mathbf{x}^k) + \sum_{i=1}^p \mu_i^k \partial_{\zeta_k} c_i(\mathbf{x}^k) \right) \quad \text{and} \quad i \notin \mathcal{I}_{\zeta_k}(\mathbf{x}^k) \Rightarrow \mu_i^k = 0. \quad (9)$$

Remembering the inclusions presented in (3) and that $\partial_{\epsilon_1} g(\mathbf{x}) \subset \partial_{\epsilon_2} g(\mathbf{x})$ for any nonsmooth function g and $0 \leq \epsilon_1 \leq \epsilon_2$, it yields $\mathbf{v}^k \in \left(\mathcal{G}_{\epsilon_k}^f(\mathbf{x}^k) + \sum_{i=1}^p \mu_i^k \mathcal{G}_{\epsilon_k}^{c_i}(\mathbf{x}^k) \right)$ and $i \notin \mathcal{I}_{\epsilon_k}(\mathbf{x}^k) \Rightarrow \mu_i^k = 0$, which proves the statement. \square

We are finally able to introduce the theorem that guarantees that the weak ϵ -ANOC is a necessary optimality condition.

Theorem 2.1. *Let \mathbf{x}^* be a local minimizer of (P). Then, \mathbf{x}^* satisfies the weak ϵ -ANOC.*

Proof. Since \mathbf{x}^* is a local minimizer of (P), there must exist $\delta_1 > 0$ such that $f(\mathbf{x}) \geq f(\mathbf{x}^*)$ for all $\mathbf{x} \in \mathcal{B}(\mathbf{x}^*, \delta_1) \cap \mathcal{F}$, where \mathcal{F} is the feasible set of (P). Moreover, setting ρ as any real number satisfying $\rho > f(\mathbf{x}^*)$, it yields that there exists $\delta_2 > 0$ with $\rho > f(\mathbf{x})$ for all $\mathbf{x} \in \mathcal{B}(\mathbf{x}^*, \delta_2)$. Choosing $\delta = \min\{\delta_1, \delta_2\}/2$, we define $\Phi_{\xi, \rho} : \mathbb{R}^n \rightarrow \mathbb{R}$ as

$$\begin{aligned} \Phi_{\xi, \rho}(\mathbf{x}) := \max \left\{ 1 - \frac{\max\{\|\mathbf{x} - \mathbf{x}^*\| - \delta, 0\}}{\xi}, 0 \right\} \max \left\{ 1 - \frac{\|\mathbf{c}(\mathbf{x})_+\|_1}{\xi}, 0 \right\} [f(\mathbf{x}) - \rho] \\ + \|\mathbf{c}(\mathbf{x})_+\|_1 + \frac{1}{2} \|\mathbf{x} - \mathbf{x}^*\|^2. \end{aligned}$$

Since $\Phi_{\xi, \rho}$ is a coercive function for any $\xi > 0$, a global minimizer of this function always exists. So, taking $\{\xi_k\} \subset \mathbb{R}_+^*$ as any sequence satisfying $\xi_k \downarrow 0$, we consider an infinite sequence $\{\mathbf{x}^k\} \subset \mathbb{R}^n$ such that $\mathbf{x}^k \in \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \Phi_{\xi_k, \rho}(\mathbf{x})$.

Consequently, $\Phi_{\xi_k, \rho}(\mathbf{x}^k) \leq \Phi_{\xi_k, \rho}(\mathbf{x}^*) = f(\mathbf{x}^*) - \rho < 0$, which yields

$$\|\mathbf{c}(\mathbf{x}^k)_+\|_1 < \xi_k \quad \text{and} \quad \|\mathbf{x}^k - \mathbf{x}^*\| < \xi_k + \delta. \quad (10)$$

This implies that the sequence $\{\mathbf{x}^k\}$ must be bounded. So, let $\hat{\mathbf{x}} \in \mathbb{R}^n$ be a cluster point of $\{\mathbf{x}^k\}$, which means that there exists an infinite index set $\mathcal{K} \subset \mathbb{N}$ such that $\mathbf{x}^k \xrightarrow[k \in \mathcal{K}]{} \hat{\mathbf{x}}$. We then define

$$0 \leq \tau_k := \max \left\{ 1 - \frac{\max\{\|\mathbf{x}^k - \mathbf{x}^*\| - \delta, 0\}}{\xi_k}, 0 \right\} \max \left\{ 1 - \frac{\|\mathbf{c}(\mathbf{x}^k)_+\|_1}{\xi_k}, 0 \right\} \leq 1.$$

Therefore, it is possible to find an infinite index set $\tilde{\mathcal{K}} \subset \mathcal{K}$ such that $\tau_k \rightarrow \hat{\tau}$, for some $\hat{\tau} \in [0, 1]$. Hence, since $\|\mathbf{c}(\hat{\mathbf{x}})_+\|_1 = 0$ because of (10), we have

$$\begin{aligned} \Phi_{\xi_k, \rho}(\mathbf{x}^k) \leq \Phi_{\xi_k, \rho}(\mathbf{x}^*) &\Rightarrow \lim_{k \in \tilde{\mathcal{K}}} \Phi_{\xi_k, \rho}(\mathbf{x}^k) \leq \lim_{k \in \tilde{\mathcal{K}}} \Phi_{\xi_k, \rho}(\mathbf{x}^*) \\ &\Rightarrow \hat{\tau} [f(\hat{\mathbf{x}}) - \rho] + \frac{1}{2} \|\hat{\mathbf{x}} - \mathbf{x}^*\|^2 \leq f(\mathbf{x}^*) - \rho. \end{aligned}$$

By (10), we know that $\|\hat{\mathbf{x}} - \mathbf{x}^*\| \leq \delta$. Then, because of the way we have defined δ , we see that $f(\hat{\mathbf{x}}) - \rho$ and $f(\mathbf{x}^*) - \rho$ are strictly negative numbers. Consequently, recalling that \mathbf{x}^* is a local minimizer of (P), we must have $\hat{\tau} = 1$ and $\hat{\mathbf{x}} = \mathbf{x}^*$. Therefore, since $\hat{\mathbf{x}}$ and $\hat{\tau}$ are arbitrary cluster points of $\{\mathbf{x}^k\}$ and $\{\tau_k\}$, respectively, this means that any cluster point of $\{\mathbf{x}^k\}$ must be \mathbf{x}^* and that any cluster point of $\{\tau_k\}$ must be 1. Since both sequences are bounded, it follows that $\mathbf{x}^k \rightarrow \mathbf{x}^*$ and $\tau_k \rightarrow 1$. In particular, this last limit also ensures that $\lim_{k \rightarrow \infty} \frac{\|\mathbf{c}(\mathbf{x}^k)_+\|_1}{\xi_k}$ exists and it is strictly less than one.

The equality $\hat{\mathbf{x}} = \mathbf{x}^*$ implies, for any sufficiently large $k \in \mathbb{N}$, that

$$\|\mathbf{x}^k - \mathbf{x}^*\| \leq \delta/2 \Rightarrow \partial \Phi_{\xi_k, \rho}(\mathbf{x}^k) = \partial \Psi_{\xi_k, \rho}(\mathbf{x}^k) + \mathbf{x}^k - \mathbf{x}^*.$$

However, by the way we have defined \mathbf{x}^k , we know that $\mathbf{0} \in \partial \Phi_{\xi_k, \rho}(\mathbf{x}^k)$, which, for any $\zeta_k \downarrow 0$, gives us

$$\mathbf{x}^* - \mathbf{x}^k \in \partial \Psi_{\xi_k, \rho}(\mathbf{x}^k) \Rightarrow \exists \mathbf{r}^k \in \partial_{\zeta_k} \Psi_{\xi_k, \rho}(\mathbf{x}^k) \text{ with } \|\mathbf{r}^k\| \rightarrow 0.$$

Due to Lemma 2.2, the statement is proven. \square

The demonstration of Lemma 2.2 gives a clue about the matter of convergence of ϵ_k . Notice that, along the proof, we had to assume $\epsilon_k/\xi_k > \zeta_k/\xi_k \rightarrow 0$, where each ξ_k is a parameter for the nonsmooth function $\Psi_{\xi_k, \rho}$. This was necessary to ensure that the multiplier associated with $\partial_{\zeta_k} f(\mathbf{x}^k)$ is bounded away from 0 in (8), which consequently means that the objective function is always considered in the sequential optimality condition (9). However, notice that even the imposition $\epsilon_k/\xi_k \rightarrow 0$ only subjects ϵ_k to the function $\Psi_{\xi_k, \rho}$, meaning that this limit is useful only

when one has a method to solve nonsmooth problems that is based on such a function, which reduces the generality of the necessary condition. This issue can be overcome by noticing that $\boldsymbol{\mu}^k$ and the limit $\epsilon_k/\xi_k \rightarrow 0$ are linked inside the proof of Lemma 2.2. So, by defining a new sequential optimality condition that requires $\epsilon_k \|\boldsymbol{\mu}^k\|_\infty \rightarrow 0$, we introduce the ϵ -Approximate Nonsmooth Optimality Condition.

Definition 2.3 (ϵ -Approximate Nonsmooth Optimality Condition). *A feasible point $\mathbf{x}^* \in \mathbb{R}^n$ of (P) is said to satisfy the ϵ -Approximate Nonsmooth Optimality Condition (ϵ -ANOC) if there exist sequences $\{\mathbf{x}^k\} \subset \mathbb{R}^n$, $\{\epsilon_k\} \subset \mathbb{R}_+$, $\{\mathbf{v}^k\} \subset \mathbb{R}^n$ and $\{\boldsymbol{\mu}^k\} \subset \mathbb{R}_+^p$ fulfilling the weak ϵ -Approximate Nonsmooth Optimality Condition at \mathbf{x}^* and, moreover, $\epsilon_k \|\boldsymbol{\mu}^k\|_\infty \rightarrow 0$ holds.*

From the proof of Theorem 2.1, one can see that the ϵ -ANOC is also a necessary optimality condition for any local solution \mathbf{x}^* of (P).

Theorem 2.2. *Let \mathbf{x}^* be a local minimizer of (P). Then, \mathbf{x}^* satisfies the ϵ -ANOC.*

Proof. Looking at the proof of Theorem 2.1, and since \mathbf{x}^* is a local minimizer of (P), one can find a sequence $\{\mathbf{x}^k\}$ converging to \mathbf{x}^* , and the associated sequences $\{\epsilon_k\}$, $\{\xi_k\}$ satisfying $\epsilon_k \downarrow 0$, $\xi_k \downarrow 0$ and $\epsilon_k/\xi_k \rightarrow 0$ such that the assumptions of Lemma 2.2 hold, and (9) is also valid (where $\zeta_k < \epsilon_k$). Additionally, since $\epsilon_k/\xi_k \rightarrow 0$, we also have $\text{conv}\{0, 1/\xi_k\}\epsilon_k \rightarrow 0$. Hence, recalling that, in (8), ι_k^f is bounded away from zero for large values of k , and that $\{\iota_k^{c_i}\epsilon_k\}$ approaches zero for all $i \in \{1, \dots, p\}$ (due to the fact that $\text{conv}\{0, 1/\xi_k\}\epsilon_k \rightarrow 0$), it must follow that $\epsilon_k \mu_i^k = \epsilon_k (\iota_k^{c_i} / \iota_k^f) \rightarrow 0$ for all $i \in \{1, \dots, p\}$. This information together with Theorem 2.1 ensure the result. \square

Remark 1. *Although the sequential optimality conditions that we have proposed so far consider only inequality constraints, one can easily generalize them to nonsmooth problems that may include equality constraints. Indeed, if in problem (P) the feasible set is given by $\mathcal{F} = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{c}(\mathbf{x}) \leq \mathbf{0}, \mathbf{h}(\mathbf{x}) = \mathbf{0}\}$, where $\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^l$, we can apply the concepts of the weak ϵ -ANOC and ϵ -ANOC to the problem*

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

where $\mathbf{g}(\mathbf{x}) = \max\{\mathbf{h}(\mathbf{x}), -\mathbf{h}(\mathbf{x})\}$. Considering (4) to the above equality-free optimization problem and using the definition of \mathbf{g} , one obtain a sequence $\{\boldsymbol{\lambda}^k\} \subset \mathbb{R}^l$ free of sign that satisfies

$$\mathbf{v}^k \in \left(\mathcal{G}_{\epsilon_k}^f(\mathbf{x}^k) + \sum_{i=1}^p \mu_i^k \mathcal{G}_{\epsilon_k}^{c_i}(\mathbf{x}^k) + \sum_{j=1}^l \lambda_j^k \mathcal{G}_{\epsilon_k}^{h_j}(\mathbf{x}^k) \right) \quad \text{and} \quad i \notin \mathcal{I}_{\epsilon_k}(\mathbf{x}^k) \Rightarrow \mu_i^k = 0.$$

For the ϵ -ANOC, one only needs to replace $\epsilon_k \|\boldsymbol{\mu}^k\|_\infty \rightarrow 0$ by $\epsilon_k \|((\boldsymbol{\mu}^k)^T, (\boldsymbol{\lambda}^k)^T)\|_\infty \rightarrow 0$.

Although the ϵ -ANOC is more restrictive over the sequence $\{\epsilon_k\}$ than the weak ϵ -ANOC, it is not clear if the former has practical advantages over the latter. The next example has the goal to show that the ϵ -ANOC is, in fact, stronger than the weak ϵ -ANOC.

Example 1. *Let us consider the following nonsmooth optimization problem*

$$\min_x f(x) = \max\{x, x^2 - 2\} \quad \text{s.t.} \quad c_1(x) = x \leq 0, c_2(x) = -x^2 \leq 0.$$

We choose $\{x^k\} \subset \mathbb{R}$ as a sequence converging to $x^* = 0$, and also let $\{\epsilon_k\}$ be such that $0 < \epsilon_k < |x^k|$. So, $2 \notin \mathcal{I}_{\epsilon_k}(x^k)$. Then, for all k sufficiently large such that $|x^k| < 1/2$, and in case $\mu_2^k = 0$, we obtain

$$\mathcal{G}_{\epsilon_k}^f(x^k) + \mu_1^k \mathcal{G}_{\epsilon_k}^{c_1}(x^k) + \mu_2^k \mathcal{G}_{\epsilon_k}^{c_2}(x^k) = 1 + \mu_1^k \geq 1.$$

This shows that if $\{x^k\}$ fulfills the required conditions of the weak ϵ -ANOC, the sequence $\{\epsilon_k\}$ must satisfy $\epsilon_k \geq |x^k|$ for every k large enough. This guarantees that $2 \in \mathcal{I}_{\epsilon_k}(x^k)$. So, for sufficiently large values of ϵ_k , one can always set $\mu_1^k = 0$, and find μ_2^k large enough in order to have $\mathcal{P}(0 | \mathcal{G}_{\epsilon_k}^f(x^k) + \mu_1^k \mathcal{G}_{\epsilon_k}^{c_1}(x^k) + \mu_2^k \mathcal{G}_{\epsilon_k}^{c_2}(x^k)) \rightarrow 0$. The previous reasoning ensures that, although x^* is a local maximum point, it satisfies the weak ϵ -ANOC.

However, we state that x^* does not satisfy the stronger condition ϵ -ANOC. By contradiction, suppose x^* satisfies the ϵ -ANOC. Then, in particular, it must fulfill the weak ϵ -ANOC, which gives us $\epsilon_k \geq |x^k|$ for every large k . Also, notice that any element in $\mathcal{G}_{\epsilon_k}^{c_2}(x^k)$ is always greater than $-2(|x^k| + \epsilon_k)$. Combining this with $\epsilon_k \geq |x^k|$ gives $g_{c_2}^k \in \mathcal{G}_{\epsilon_k}^{c_2}(x^k) \Rightarrow g_{c_2}^k \geq -4\epsilon_k$. So, $v^k \in (\mathcal{G}_{\epsilon_k}^f(x^k) + \mu_1^k \mathcal{G}_{\epsilon_k}^{c_1}(x^k) + \mu_2^k \mathcal{G}_{\epsilon_k}^{c_2}(x^k))$ implies

$$v^k = 1 + \mu_1^k + \mu_2^k g_{c_2}^k \geq 1 - 4\mu_2^k \epsilon_k. \quad (11)$$

Since we supposed that x^* satisfies the ϵ -ANOC, it follows $\epsilon_k \|\mu^k\|_\infty \rightarrow 0$. This last fact together with (11) show that v^k cannot converge to zero, which evinces a contradiction with the statement that x^* satisfies the ϵ -ANOC.

We now prove that our proposed optimality conditions must be, at least, as strong as the necessary optimality condition presented in [29] (considering $D(\mathbf{x}^*) = \mathbb{R}^n$ in [29, Proposition 3.1]), which has its roots in the study presented in [36].

Theorem 2.3. *If $\mathbf{x}^* \in \mathbb{R}^n$ is a feasible point for (P) satisfying the (weak) ϵ -ANOC, then the generalized Fritz-John optimality conditions hold at \mathbf{x}^* , i.e., there must exist positive real values $\lambda_0, \lambda_1, \dots, \lambda_p$, not all simultaneously zero, such that*

$$\mathbf{0} \in \lambda_0 \partial f(\mathbf{x}^*) + \sum_{i=1}^p \lambda_i \partial c_i(\mathbf{x}^*) \quad \text{and} \quad \lambda_i c_i(\mathbf{x}^*) = 0, \text{ for all } i \in \{1, \dots, p\}. \quad (12)$$

Proof. For any feasible point $\mathbf{x}^* \in \mathbb{R}^n$ satisfying the (weak) ϵ -ANOC, there exist sequences $\{\mathbf{x}^k\} \subset \mathbb{R}^n$, $\{\epsilon_k\} \subset \mathbb{R}_+^*$, $\{\mathbf{v}^k\} \subset \mathbb{R}^n$ and $\{\mu^k\} \subset \mathbb{R}_+^p$ such that $\mathbf{x}^k \rightarrow \mathbf{x}^*$, $\epsilon_k \downarrow 0$ and $\|\mathbf{v}^k\| \rightarrow 0$, with

$$\mathbf{v}^k = \left(\mathbf{v}_f^k + \sum_{i=1}^p \mu_i^k \mathbf{v}_{c_i}^k \right) \quad \text{and} \quad i \notin \mathcal{I}_{\epsilon_k}(\mathbf{x}^k) \Rightarrow \mu_i^k = 0, \quad (13)$$

where $\mathbf{v}_f^k \in \mathcal{G}_{\epsilon_k}^f(\mathbf{x}^k) \subset \partial_{\epsilon_k} f(\mathbf{x}^k)$ and $\mathbf{v}_{c_i}^k \in \mathcal{G}_{\epsilon_k}^{c_i}(\mathbf{x}^k) \subset \partial_{\epsilon_k} c_i(\mathbf{x}^k)$. Then two situations can happen: (a) the sequence $\{\mu^k\}$ is bounded, or (b) the sequence $\{\mu^k\}$ is unbounded.

Let us consider that (a) is the case. Since $\mathbf{x}^k \rightarrow \mathbf{x}^*$ and all functions considered here are locally Lipschitz continuous functions, it follows that $\{\mathbf{v}_f^k\}$ and $\{\mathbf{v}_{c_i}^k\}$, $i \in \{1, \dots, p\}$, are all bounded sequences. Therefore, there must exist an infinite index set $\mathcal{K} \subset \mathbb{N}$ such that $\mu^k \xrightarrow[k \in \mathcal{K}]{} \mu^*$, $\mathbf{v}_f^k \xrightarrow[k \in \mathcal{K}]{} \mathbf{v}_f^*$ and $\mathbf{v}_{c_i}^k \xrightarrow[k \in \mathcal{K}]{} \mathbf{v}_{c_i}^*$, $i \in \{1, \dots, p\}$. So considering the result of [40, Lemma 3.2 (iii)] for the auxiliary functions $\tilde{f}(\mathbf{x}) = f(\mathbf{x}) - \mathbf{v}_f^{*T} \mathbf{x}$ and $\tilde{c}_i(\mathbf{x}) = c_i(\mathbf{x}) - \mathbf{v}_{c_i}^{*T} \mathbf{x}$, $i \in \{1, \dots, p\}$, it follows that $\mathbf{v}_f^* \in \partial f(\mathbf{x}^*)$ and $\mathbf{v}_{c_i}^* \in \partial c_i(\mathbf{x}^*)$, $i \in \{1, \dots, p\}$. Therefore, recalling that $\|\mathbf{v}^k\| \rightarrow 0$, we get $\mathbf{0} \in (\partial f(\mathbf{x}^*) + \sum_{i=1}^p \mu_i^* \partial c_i(\mathbf{x}^*))$. Notice also that if $c_i(\mathbf{x}^*) < 0$ for some i , then $i \notin \mathcal{I}_{\epsilon_k}(\mathbf{x}^k) (\Rightarrow \mu_i^k = 0)$ for all $k \in \mathbb{N}$ sufficiently large, which yields $\mu_i^* = 0$. This assures the result for case (a).

If (b) holds, then one only has to divide the equation that appears in (13) by $\max_i \{\mu_i^k\}$. This ensures that the new multipliers of $\mathbf{v}_{c_i}^k$, $i \in \{1, \dots, p\}$, will be bounded and the multiplier of \mathbf{v}_f^k will be $1/\max_i \{\mu_i^k\}$. Hence, the same reasoning used in case (a) can be employed, obtaining the Fritz-John conditions with zero being the multiplier of $\partial f(\mathbf{x}^*)$. \square

Considering the Mangasarian-Fromovitz constraint qualification for nonsmooth problems [27, Section 3], which guarantees that (12) cannot be true if $\lambda_0 = 0$, and the generalized KKT condition, as defined below, the corollary presented in the sequence follows immediately from the theorem above.

Definition 2.4 (Generalized KKT conditions [18, Section 6.3]). *A feasible point $\mathbf{x}^* \in \mathbb{R}^n$ of (P) is said to satisfy the generalized KKT conditions if there exists $\mu^* \in \mathbb{R}_+^p$ such that*

$$\mathbf{0} \in \left(\partial f(\mathbf{x}^*) + \sum_{i=1}^p \mu_i^* \partial c_i(\mathbf{x}^*) \right) \quad \text{and} \quad \mu_i^* c_i(\mathbf{x}^*) = 0 \text{ for all } i \in \{1, \dots, p\}.$$

Corollary 2.1. *If a feasible point $\mathbf{x}^* \in \mathbb{R}^n$ for (P) satisfies the (weak) ϵ -ANOC, and the nonsmooth Mangasarian-Fromovitz constraint qualification holds at \mathbf{x}^* , then the generalized KKT condition also holds at \mathbf{x}^* .*

We believe that the results obtained in this section have great relevance in the nonsmooth optimization field, specially when one is concerned with convergence analysis of nonsmooth optimization methods. In the nonsmooth context, it is common to prove the convergence of an algorithm by showing that every accumulation point generated by the iterates of such a method is a stationary point of (Unc-P) (see, for example, [23, 38, 39]). This can be well justified when a property called *calmness* holds at a local minimizer \mathbf{x}^* of the original problem (P).

Definition 2.5 (Calmness [12, 32]). *Consider \mathbf{x}^* a local minimizer of (P). We say that problem (P) is calm at \mathbf{x}^* when one can find $\rho > 0$ sufficiently large such that \mathbf{x}^* is also a locally optimal solution of (Unc-P).*

Thus, under the calmness hypothesis, the statement that \mathbf{x}^* is a locally optimal solution of (Unc-P) becomes a necessary optimality condition. However, this is not a legitimate necessary optimality condition, in the sense that one needs to assume more than just the fact that \mathbf{x}^* is a local minimizer of (P). That is not the case for the (weak) ϵ -ANOC. By the results of Theorems 2.1 and 2.2, we have shown that the (weak) ϵ -ANOC at \mathbf{x}^* is a legitimate necessary optimality condition. Therefore, the results exposed here give alternative forms of dealing with the convergence analysis of nonsmooth optimization methods even in the absence of calmness.

We end our theoretical analysis with a brief section presenting the relations between the ϵ -ANOC and the sequential optimality conditions AKKT and CAKKT [2, 4] in the smooth context.

3 ϵ -ANOC applied to smooth optimization

The goal of this section is to show that, under mild assumptions, the scheme from Fig. 3 holds true when the objective and constraints functions of (P) are all differentiable. For completeness, we define below the concepts of AKKT and CAKKT for the following smooth optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \quad \text{s.t.} \quad \mathbf{c}(\mathbf{x}) \leq \mathbf{0}, \quad \mathbf{h}(\mathbf{x}) = \mathbf{0}, \quad (\text{P-smooth})$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $\mathbf{c} : \mathbb{R}^n \rightarrow \mathbb{R}^p$ and $\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^l$ are all functions of class C^1 .

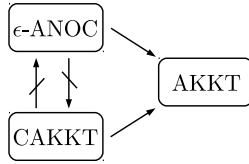


Figure 2: Scheme of relations between the ϵ -ANOC and the well-known sequential optimality conditions in the smooth field.

Definition 3.1 (AKKT, CAKKT [2, 4]). *A feasible point \mathbf{x}^* of (P-smooth) is an AKKT point if there exist sequences $\{\mathbf{x}^k\}$, $\{\boldsymbol{\mu}^k\} \subset \mathbb{R}_+^p$ and $\{\boldsymbol{\lambda}^k\} \subset \mathbb{R}^l$ such that $\mathbf{x}^k \rightarrow \mathbf{x}^*$, $\lim_{k \rightarrow \infty} \min\{-c_i(\mathbf{x}^k), \mu_i^k\} = 0$ for all $i \in \{1, \dots, p\}$, and*

$$\lim_{k \rightarrow \infty} \left\| \nabla f(\mathbf{x}^k) + \sum_{i=1}^p \mu_i^k \nabla c_i(\mathbf{x}^k) + \sum_{j=1}^l \lambda_j^k \nabla h_j(\mathbf{x}^k) \right\| = 0.$$

If, in addition, $\lim_{k \rightarrow \infty} \mu_i^k c_i(\mathbf{x}^k) = 0$ and $\lim_{k \rightarrow \infty} \lambda_j^k h_j(\mathbf{x}^k) = 0$ for all $i \in \{1, \dots, p\}$ and $j \in \{1, \dots, l\}$, then we call \mathbf{x}^ a CAKKT point.*

First, we prove that every feasible point \mathbf{x}^* of any smooth optimization problem satisfying the ϵ -ANOC must also be an AKKT point.

Theorem 3.1. *Suppose the functions \mathbf{c} and \mathbf{h} in (P-smooth) have locally Lipschitz continuous derivatives. If x^* is a feasible point of (P-smooth) and satisfies the ϵ -ANOC, then x_* is also an AKKT point.*

Proof. Since x^* satisfies the ϵ -ANOC, there must exist sequences $\{\mathbf{x}^k\} \subset \mathbb{R}^n$, $\{\epsilon_k\} \subset \mathbb{R}_+^*$, $\{\mathbf{v}^k\} \subset \mathbb{R}^n$, $\{\boldsymbol{\mu}^k\} \subset \mathbb{R}_+^p$ and $\{\boldsymbol{\lambda}^k\} \subset \mathbb{R}^l$ such that $\mathbf{x}^k \rightarrow \mathbf{x}^*$, $\epsilon_k \downarrow 0$ and $\|\mathbf{v}^k\| \rightarrow 0$, where

$$\epsilon_k \|((\boldsymbol{\mu}^k)^T, (\boldsymbol{\lambda}^k)^T)\|_\infty \rightarrow 0,$$

$i \notin \mathcal{I}_{\epsilon_k}(\mathbf{x}^k) \Rightarrow \mu_i^k = 0$ and

$$\mathbf{v}^k \in \left(\mathcal{G}_{\epsilon_k}^f(\mathbf{x}^k) + \sum_{i=1}^p \mu_i^k \mathcal{G}_{\epsilon_k}^{c_i}(\mathbf{x}^k) + \sum_{j=1}^l \lambda_j^k \mathcal{G}_{\epsilon_k}^{h_j}(\mathbf{x}^k) \right). \quad (14)$$

Since the derivatives of the involved functions are locally Lipschitz continuous, it follows that

$$\mathbf{v}^k = \nabla f(\mathbf{x}^k) + \sum_{i=1}^p \mu_i^k \nabla c_i(\mathbf{x}^k) + \sum_{j=1}^l \lambda_j^k \nabla h_j(\mathbf{x}^k) + \mathbf{r}_f^k + \sum_{i=1}^p \mu_i^k \mathbf{r}_{c_i}^k + \sum_{j=1}^l \lambda_j^k \mathbf{r}_{h_j}^k,$$

where \mathbf{r}_f^k , $\mathbf{r}_{c_i}^k$ and $\mathbf{r}_{h_j}^k$ are all error vectors of order ϵ_k . Since $\|\mathbf{v}^k\| \rightarrow 0$, $\epsilon_k \downarrow 0$ and $\epsilon_k \|((\boldsymbol{\mu}^k)^T, (\boldsymbol{\lambda}^k)^T)\|_\infty \rightarrow 0$, we get

$$\nabla f(\mathbf{x}^k) + \sum_{i=1}^p \mu_i^k \nabla c_i(\mathbf{x}^k) + \sum_{j=1}^l \lambda_j^k \nabla h_j(\mathbf{x}^k) \rightarrow \mathbf{0}. \quad (15)$$

Additionally, because $i \notin \mathcal{I}_{\epsilon_k}(\mathbf{x}^k) \Rightarrow \mu_i^k = 0$, one can see that $\min\{-c_i(\mathbf{x}^k), \mu_i^k\} \rightarrow 0$ for all $i \in \{1, \dots, p\}$. This together with (15) fulfill the conditions of an AKKT point, which completes the proof. \square

Remark 2. *It is worth noticing that the reciprocal is not true, i.e., AKKT does not imply the ϵ -ANOC in the smooth case. Indeed, considering the optimization problem given in Example 1, $x^* = 0$ is an AKKT point (around x^* the optimization problem is smooth, that justifies applying the AKKT concept at this point). However, we have showed that x^* does not satisfy the ϵ -ANOC, which guarantees that $\text{AKKT} \not\Rightarrow \epsilon\text{-ANOC}$. Hence, our new sequential optimality condition is also stronger than AKKT in the smooth context.*

We now proceed with two examples showing that CAKKT and the ϵ -ANOC are not connected to each other.

Example 2. *Let us consider the following constrained optimization problem*

$$\min_{\mathbf{x}} f(\mathbf{x}) = x_1 - x_2 \quad \text{s.t.} \quad x_1 \geq 0, \quad x_2 \geq 0, \quad x_1 x_2 \leq 0.$$

Although the point $\mathbf{x}^* = \mathbf{0}$ is not a local minimizer for the problem, if one chooses $\mathbf{x}^k = (1/k, -1/k)^T$, $\mu_1^k = \mu_2^k = 0$ and $\mu_3^k = k$, one can see that \mathbf{x}^* is a CAKKT point [3]. However, we state that \mathbf{x}^* does not fulfill the ϵ -ANOC.

Indeed, suppose, by contradiction, that \mathbf{x}^* satisfies the ϵ -ANOC. Recalling the arguments used to go from (14) to (15), and since the objective and constraint functions are all smooth and have locally Lipschitz continuous derivatives, there must exist $\{\mathbf{x}^k\} = \{(x_1^k, x_2^k)^T\}$ and $\{\boldsymbol{\mu}^k\} \subset \mathbb{R}_+^3$ with $\epsilon_k \|\boldsymbol{\mu}^k\|_\infty \rightarrow 0$ such that $\mathbf{x}^k \rightarrow \mathbf{x}^*$ and

$$\begin{bmatrix} 1 \\ -1 \end{bmatrix} + \mu_1^k \begin{bmatrix} -1 \\ 0 \end{bmatrix} + \mu_2^k \begin{bmatrix} 0 \\ -1 \end{bmatrix} + \mu_3^k \begin{bmatrix} x_2^k \\ x_1^k \end{bmatrix} \rightarrow \mathbf{0}.$$

By the second line of the above relation, we have that, for large values of k , the following must hold: $x_1^k > 0$ and $\mu_3^k > 1/(2x_1^k)$. Since $\epsilon_k \mu_3^k \rightarrow 0$, it yields $\epsilon_k = o(x_1^k)$. Hence, $1 \notin \mathcal{I}_{\epsilon_k}(\mathbf{x}^k)$ and, consequently, $\mu_1^k = 0$ for large values of k . So, by the first line, we obtain $\mu_3^k x_2^k \rightarrow -1$, which assures $x_2^k < 0$ and $\mu_3^k > 1/(2|x_2^k|)$ for k sufficiently large. Since μ_3^k can only be different from zero if $3 \in \mathcal{I}_{\epsilon_k}(\mathbf{x}^k)$, the value ϵ_k must be greater than $\min\{x_1^k, |x_2^k|\}$ for large values of k . This contradicts the fact that $\epsilon_k \|\boldsymbol{\mu}^k\|_\infty \rightarrow 0$, which allows us to conclude that \mathbf{x}^* does not satisfy the ϵ -ANOC.

Example 3. We now consider the smooth optimization problem

$$\min_{\mathbf{x}} f(\mathbf{x}) = \frac{(x_2 - 2)^2}{2} \quad \text{s.t.} \quad x_1 = 0, \quad x_1 x_2 = 0.$$

As shown in [4], the point $\mathbf{x}^* = (0, 1)^T$ is not a CAKKT point. However, if we choose $\mathbf{x}^k = (1/k, 1)^T$, $\boldsymbol{\lambda} = (-k, k)^T$ and $\epsilon_k = 1/k^2$, we obtain that the ϵ -ANOC is fulfilled as defined in Remark 1. Therefore, \mathbf{x}^* satisfies the ϵ -ANOC.

The above examples show that, in fact, the ϵ -ANOC is a new legitimate (i.e., that does not require any kind of constraint qualification) necessary optimality condition even in the case where all functions are smooth.

4 A practical algorithm

The goal of this section is to present an algorithm that is able to generate a sequence of iterates satisfying the sequential optimality conditions introduced in the last section. Looking at relation (4) and the proof of Theorem 2.1, one can easily come up with the following algorithm: given $\xi_k \downarrow 0$ and $\rho \in \mathbb{R}$ sufficiently large, whenever possible, set

$$\mathbf{x}^k \in \underset{\mathbf{x} \in \mathbb{R}^n}{\operatorname{argmin}} \Psi_{\xi_k, \rho}(\mathbf{x}). \quad (16)$$

Although conceptually simple, the above procedure cannot be used in a practical method. For optimization problems involving nonconvex functions, an implementable algorithm usually guarantees that the method will only find stationary points for the problem, not minimizers. Moreover, the set of minimizers of the function $\Psi_{\xi_k, \rho}$ may be empty, therefore this procedure may not be well defined.

We thus present Algorithm 1 (PACNO - Penalized Algorithm for Constrained Nonsmooth Optimization), where we avoid these issues by using a solution of a relaxed version of (16) (see Remark 3 for a more detailed discussion about this algorithm) where we trade the global minimization for the easier problem of Step 1 of Algorithm 1, which concerns only with approximate stationary points. This is sufficient to ensure that Algorithm 1 generates a sequence of iterates that verifies the ϵ -ANOC. However, first we need to show that the set \mathcal{S}_{Ψ}^k defined in (17) will eventually be nonempty and, moreover, that the sequence $\{\rho_k\}$ generated by the algorithm will stabilize at a sufficiently large value. To guarantee that $\mathcal{S}_{\Psi}^k \neq \emptyset$, we need to consider one extra assumption.

Assumption 1. *There exists $\alpha > 0$ such that $\mathcal{F}_\alpha := \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{c}(\mathbf{x})_+\|_1 \leq \alpha\}$ is compact.*

Since we are dealing with possibly nonconvex functions, it is not sufficient to ask for the compactness of the feasible set, but, instead, one needs to ask for the compactness of the perturbed feasible set \mathcal{F}_α . This is required to exclude some pathological functions like the one illustrated in Fig. 3, where the feasible set is compact but there is no $\alpha > 0$ ensuring the compactness of \mathcal{F}_α . Nevertheless, the assumption can be easily ensured if the feasible points are additionally restricted to box constraints. Noticing that unbounded solutions do not occur for the majority of real problems, the user may include sufficiently large artificial box constraints to enforce the condition.

We start by showing a technical lemma that will be helpful in the subsequent results.

Algorithm 1: Penalized Algorithm for Constrained Nonsmooth Optimization

Step 0. Choose $\mathbf{x}^0 \in \mathbb{R}^n$, $\rho_1 \in [0, +\infty)$, $M \in (0, +\infty)$, $\{\theta_\xi, \theta_\epsilon, \theta_\nu\} \subset (0, 1)$, $\omega \in (0, 1]$, $\{\xi_1, \epsilon_1, \nu_1\} \subset (0, +\infty)$, $\xi_{\text{opt}} \in [0, \xi_1]$, $\epsilon_{\text{opt}} \in [0, \epsilon_1]$ and $\nu_{\text{opt}} \in [0, \nu_1]$.
Set $k = 1$.

Step 1. Define

$$\mathcal{S}_\Psi^k := \{\mathbf{x} : \Psi_{\xi_k, \rho_k}(\mathbf{x}) \leq \Psi_{\xi_k, \rho_k}(\mathbf{x}^{k-1}) \quad \text{and} \quad \|\mathcal{P}(\mathbf{0} \mid \partial_{\epsilon_k} \Psi_{\xi_k, \rho_k}(\mathbf{x}))\| \leq \nu_k\}. \quad (17)$$

If $\mathcal{S}_\Psi^k \neq \emptyset$, then set \mathbf{x}^k as any point in \mathcal{S}_Ψ^k . Otherwise, set \mathbf{x}^k as any point \mathbf{x} satisfying $\|\mathbf{c}(\mathbf{x})_+\|_1 \leq \|\mathbf{c}(\mathbf{x}^{k-1})_+\|_1$ and $\|\mathcal{P}(\mathbf{0} \mid \partial_{\epsilon_k} \|\mathbf{c}(\mathbf{x})_+\|_1)\| \leq \nu_k$, and go to Step 3.

Step 2. If $\|\mathbf{c}(\mathbf{x}^k)_+\|_1 < \xi_{\text{opt}}$, then

$$\begin{cases} \text{terminate} & , \text{ if } \epsilon_k \leq \epsilon_{\text{opt}} \text{ and } \nu_k \leq \nu_{\text{opt}} \\ \text{set } \xi_{k+1} = \xi_k \text{ and go to Step 5} & , \text{ otherwise.} \end{cases}$$

Step 3. If $\|\mathbf{c}(\mathbf{x}^k)_+\|_1 > \omega \xi_k$ with $\epsilon_k \leq \epsilon_{\text{opt}}$ and $\nu_k \leq \nu_{\text{opt}}$, then stop. The iterate \mathbf{x}^k is close to a stationary point of the infeasibility measure.

Step 4. Set $\xi_{k+1} = \xi_k - (1 - \theta_\xi)[\xi_k - \|\mathbf{c}(\mathbf{x}^k)_+\|_1]_+$.

Step 5. If $f(\mathbf{x}^k) - \rho_k \leq -M$, then set $\rho_{k+1} = \rho_k$. Otherwise, proceed with $\rho_{k+1} = \rho_k + 2(M + [f(\mathbf{x}^k) - \rho_k]_+)$.

Step 6. Set $\epsilon_{k+1} = \theta_\epsilon \epsilon_k$, $\nu_{k+1} = \theta_\nu \nu_k$ and $k \leftarrow k + 1$. Go back to Step 1.

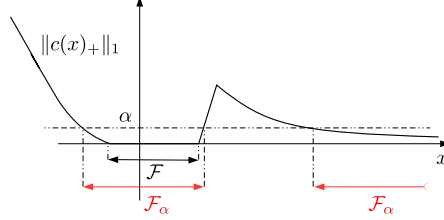


Figure 3: An illustration of a function that does not have a compact \mathcal{F}_α . The set \mathcal{F} stands for the feasible region. Notice that one cannot find $\alpha > 0$ that ensures the compactness of \mathcal{F}_α .

Lemma 4.1. Let $\{\mathbf{x}^k\}$ be a sequence of iterates generated by Algorithm 1 with $\xi_{\text{opt}} = \epsilon_{\text{opt}} = \nu_{\text{opt}} = 0$. Then, there must exist a sequence $\{\delta_k\} \subset \mathbb{R}_+^*$ converging to zero such that $\xi_k \leq \|\mathbf{c}(\mathbf{x}^k)_+\|_1 + \delta_k$. Moreover, if $\liminf_{k \rightarrow \infty} \|\mathbf{c}(\mathbf{x}^k)_+\|_1 = 0$, then $\xi_k \downarrow 0$.

Proof. Since $\xi_{\text{opt}} = \epsilon_{\text{opt}} = \nu_{\text{opt}} = 0$, the algorithm will have infinitely many iterations. Moreover, notice that $\{\xi_k\}$ is a positive bounded monotone decreasing sequence, which ensures that such a sequence must converge to some positive real value $\bar{\xi}$. Hence, since ξ_k is always updated by Step 4 (when $\xi_{\text{opt}} = 0$), we have

$$\begin{aligned} \xi_{k+1} = \xi_k - (1 - \theta_\xi)[\xi_k - \|\mathbf{c}(\mathbf{x}^k)_+\|_1]_+ &\Rightarrow (1 - \theta_\xi)[\xi_k - \|\mathbf{c}(\mathbf{x}^k)_+\|_1]_+ = \xi_k - \xi_{k+1} \\ &\Rightarrow \lim_{k \rightarrow \infty} [\xi_k - \|\mathbf{c}(\mathbf{x}^k)_+\|_1]_+ = 0. \end{aligned}$$

In other words, there exists $\{\delta_k\} \subset \mathbb{R}_+^*$, with $\delta_k \rightarrow 0$, such that $\xi_k \leq \|\mathbf{c}(\mathbf{x}^k)_+\|_1 + \delta_k$. So, if $\liminf_{k \rightarrow \infty} \|\mathbf{c}(\mathbf{x}^k)_+\|_1 = 0$, it is straightforward to see that $\xi_k \downarrow 0$. \square

The next lemma is useful to prove the result that, eventually, the set \mathcal{S}_Ψ^k will have at least one element. It gives a sufficient condition for the existence of a global minimizer of $\Psi_{\xi, \rho}$.

Lemma 4.2. *Let ρ be a real number such that*

$$\min_{\mathbf{x} \in \mathcal{F}_{\alpha/2}} f(\mathbf{x}) - \rho < 0$$

and $\xi \in (0, \alpha]$, with α being a strictly positive real number satisfying Assumption 1. Then, the function $\Psi_{\xi, \rho}$ defined in (5) has a global minimizer.

Proof. Consider any fixed $\xi \in (0, \alpha]$. Let us first show that $\Psi_{\xi, \rho}$ is bounded from below. By contradiction, let us suppose the opposite, i.e., we can find a sequence $\{\mathbf{x}^k\}$ such that $\Psi_{\xi, \rho}(\mathbf{x}^k) \leq -k$. Since $-k$ is a negative number, we must have, for all $k \in \mathbb{N}$, that $\|\mathbf{c}(\mathbf{x}^k)_+\|_1 \leq \xi$, since, otherwise, we would have $\Psi_{\xi, \rho}(\mathbf{x}^k) = \|\mathbf{c}(\mathbf{x}^k)_+\|_1 \geq 0$. Because $\xi \in (0, \alpha]$, this means that $\mathbf{x}^k \in \mathcal{F}_\alpha$, for all $k \in \mathbb{N}$. However, $\Psi_{\xi, \rho}$ is a continuous function, which yields that it must assume a minimum value on \mathcal{F}_α . This is a contradiction with the fact that $\Psi_{\xi, \rho}(\mathbf{x}^k) \leq -k$, for all $k \in \mathbb{N}$. Therefore, the function $\Psi_{\xi, \rho}$ must be bounded.

Now, let us then consider a sequence $\{\mathbf{z}^k\}$ such that $\Psi_{\xi, \rho}(\mathbf{z}^k) \rightarrow \inf_{\mathbf{x} \in \mathbb{R}^n} \Psi_{\xi, \rho}(\mathbf{x})$. By hypothesis, we know that there exists $\hat{\mathbf{x}} \in \mathcal{F}_{\alpha/2}$ such that $f(\hat{\mathbf{x}}) - \rho < 0$. This implies that

$$\inf_{\mathbf{x} \in \mathbb{R}^n} \Psi_{\xi, \rho}(\mathbf{x}) \leq \Psi_{\xi, \rho}(\hat{\mathbf{x}}) \leq \|\mathbf{c}(\hat{\mathbf{x}})_+\|_1.$$

Therefore, for any large $k \in \mathbb{N}$, we must have $\Psi_{\xi, \rho}(\mathbf{z}^k) \leq 2\|\mathbf{c}(\hat{\mathbf{x}})_+\|_1$. Since, $\hat{\mathbf{x}} \in \mathcal{F}_{\alpha/2}$, it implies that $\Psi_{\xi, \rho}(\mathbf{z}^k) \leq \alpha$. Recalling that $\xi \leq \alpha$, we must have $\mathbf{z}^k \in \mathcal{F}_\alpha$, for any large $k \in \mathbb{N}$. So, because \mathcal{F}_α is compact, one can ensure that there exists a subsequence of $\{\mathbf{z}^k\}$ converging to $\bar{\mathbf{z}} \in \mathcal{F}_\alpha$ such that $\Psi_{\xi, \rho}(\bar{\mathbf{z}}) = \inf_{\mathbf{x} \in \mathbb{R}^n} \Psi_{\xi, \rho}(\mathbf{x})$, which ends the proof. \square

The following result is a technical lemma and it proves that $\|\mathbf{c}(\mathbf{x}^k)_+\|_1$ will be related with the value ξ_k and the infeasibility measure of the past iteration. It also ensures that our method will not suffer from the greediness phenomenon.

Lemma 4.3. *Let $k \in \mathbb{N}$ be any iteration of Algorithm 1. Then,*

$$\|\mathbf{c}(\mathbf{x}^k)_+\|_1 \leq \max \{ \xi_k, \|\mathbf{c}(\mathbf{x}^{k-1})_+\|_1 \}.$$

Proof. By contradiction, suppose that there exists an iteration k such that

$$\|\mathbf{c}(\mathbf{x}^k)_+\|_1 > \max \{ \xi_k, \|\mathbf{c}(\mathbf{x}^{k-1})_+\|_1 \}. \quad (18)$$

Then, because of Step 1, we have $\mathbf{x}^k \in \mathcal{S}_\Psi^k$. Therefore, $\Psi_{\xi_k, \rho_k}(\mathbf{x}^k) \leq \Psi_{\xi_k, \rho_k}(\mathbf{x}^{k-1})$. Moreover, because of Step 5, one can see that $f(\mathbf{x}^{k-1}) - \rho_k < 0$, which implies

$$\Psi_{\xi_k, \rho_k}(\mathbf{x}^k) \leq \Psi_{\xi_k, \rho_k}(\mathbf{x}^{k-1}) \leq \|\mathbf{c}(\mathbf{x}^{k-1})_+\|_1.$$

Thus, by (18), we have $\Psi_{\xi_k, \rho_k}(\mathbf{x}^k) < \|\mathbf{c}(\mathbf{x}^k)_+\|_1$. Consequently, the above inequality can only be true in the case that $\|\mathbf{c}(\mathbf{x}^k)_+\|_1 < \xi_k$, which contradicts (18). \square

The next result tells us that, when the sequence of iterates approaches the feasible set, then, eventually, the set \mathcal{S}_Ψ^k will be nonempty and the sequence $\{\rho_k\}$ will stabilize.

Lemma 4.4. *Let $\{\mathbf{x}^k\}$ be a sequence of iterates generated by Algorithm 1 with $\xi_{opt} = \epsilon_{opt} = \nu_{opt} = 0$. If $\liminf_{k \rightarrow \infty} \|\mathbf{c}(\mathbf{x}^k)_+\|_1 = 0$, then there exist a sufficiently large $\hat{k} \in \mathbb{N}$ and $\bar{\rho} \in \mathbb{R}$ such that, for all $k \geq \hat{k}$, we have*

$$\mathcal{S}_\Psi^k \neq \emptyset \quad \text{and} \quad \rho_k = \bar{\rho}.$$

Proof. First, let us prove that $\mathcal{S}_\Psi^k \neq \emptyset$. Since, $\liminf_{k \rightarrow \infty} \|\mathbf{c}(\mathbf{x}^k)_+\|_1 = 0$, we must have an iteration \hat{k} such that $\|\mathbf{c}(\mathbf{x}^{\hat{k}})_+\|_1 \leq \alpha/2$, where α is the positive real number satisfying Assumption 1. Consequently, by Step 5 of the algorithm, one can see that

$$\min_{\mathbf{x} \in \mathcal{F}_{\alpha/2}} f(\mathbf{x}) - \rho_k \leq f(\mathbf{x}^{\hat{k}}) - \rho_k < 0, \text{ for any } k > \hat{k}.$$

Additionally, Lemma 4.1 tells us that $\xi_k \rightarrow 0$. So, by Lemma 4.2, we see, for any large $k \in \mathbb{N}$, that $\mathcal{S}_\Psi^k \neq \emptyset$.

Let us now prove that there exists $\bar{\rho} \in \mathbb{R}$ such that $\rho_k = \bar{\rho}$, for any large $k \in \mathbb{N}$. By contradiction, let us assume that the statement is false. This means, by the way the algorithm was designed, that there exists an infinite index set $\mathcal{K} \subset \mathbb{N}$ such that

$$f(\mathbf{x}^k) - \rho_k > -M, \text{ for all } k \in \mathcal{K}, \quad (19)$$

implying that $\rho_k \rightarrow \infty$. Because $\liminf_{k \rightarrow \infty} \|\mathbf{c}(\mathbf{x}^k)_+\|_1 = 0$, it means that there exists an iteration \hat{k} satisfying $\|\mathbf{c}(\mathbf{x}^{\hat{k}})_+\|_1 \leq \alpha$ and $\xi_{\hat{k}} \leq \alpha$ (because of Lemma 4.1). Consequently, due to Lemma 4.3 and recalling that $\{\xi_k\}$ is a monotone decreasing sequence, we have

$$\|\mathbf{c}(\mathbf{x}^k)_+\|_1 \leq \alpha, \text{ i.e., } \mathbf{x}^k \in \mathcal{F}_\alpha, \text{ for all } k \geq \hat{k}.$$

However, by hypothesis, \mathcal{F}_α is compact, which, together with the assumption $\rho_k \rightarrow \infty$, implies

$$\rho_k \geq \max_{\mathbf{x} \in \mathcal{F}_\alpha} f(\mathbf{x}) + M \geq f(\mathbf{x}^k) + M, \text{ for any large } k \in \mathbb{N}. \quad (20)$$

This is a contradiction with (19). Therefore, there must exist $\bar{\rho}$ such that $\rho_k = \bar{\rho}$, for any large $k \in \mathbb{N}$. \square

We are now ready to present the convergence theorem of Algorithm 1. Because the nonsmooth optimization problem in hand might involve nonconvex functions, the result ensures that a cluster point \mathbf{x}^* of the iterate sequence $\{\mathbf{x}^k\}$ can either be a stationary point for the infeasibility measure $\|\mathbf{c}(\cdot)_+\|_1$, or satisfies the ϵ -ANOC.

Theorem 4.1. *If $\{\mathbf{x}^k\}$ is the sequence of iterates generated by Algorithm 1 with $\xi_{\text{opt}} = \epsilon_{\text{opt}} = \nu_{\text{opt}} = 0$ and $\theta_\epsilon \in (0, \theta_\xi)$, then, given any infinite index set $\mathcal{K} \subset \mathbb{N}$ such that $\mathbf{x}^k \xrightarrow[k \in \mathcal{K}]{} \mathbf{x}^*$, for some $\mathbf{x}^* \in \mathbb{R}^n$, one of the following statements must be true:*

- a) \mathbf{x}^* satisfies $0 \in \partial\|\mathbf{c}(\mathbf{x}^*)_+\|_1$;
- b) \mathbf{x}^* is a feasible point for (P) and it satisfies the ϵ -ANOC. Moreover, there exists $\bar{\rho} \in \mathbb{R}$ such that $\rho_k = \bar{\rho}$, for any large $k \in \mathbb{N}$.

Proof. We start by noticing that, since $\xi_{\text{opt}} = \epsilon_{\text{opt}} = \nu_{\text{opt}} = 0$, Algorithm 1 in fact generates an infinite sequence of points $\{\mathbf{x}^k\} \subset \mathbb{R}^n$. Moreover, ξ_k is always updated by Step 4 (since $\xi_{\text{opt}} = 0$), hence $\xi_{k+1} = \xi_k - (1 - \theta_\xi)[\xi_k - \|\mathbf{c}(\mathbf{x}^k)_+\|_1]_+ \geq \theta_\xi \xi_k$. Therefore, because we suppose $\theta_\epsilon \in (0, \theta_\xi)$, it yields

$$\frac{\epsilon_{k+1}}{\xi_{k+1}} \leq \eta \frac{\epsilon_k}{\xi_k} \leq \eta^2 \frac{\epsilon_{k-1}}{\xi_{k-1}} \leq \dots \leq \eta^k \frac{\epsilon_1}{\xi_1}, \text{ where } \eta = \frac{\theta_\epsilon}{\theta_\xi} < 1,$$

which ensures $\epsilon_k/\xi_k \rightarrow 0$. We then divide the proof in the following cases:

- i) $\liminf_{k \in \mathcal{K}} \|\mathbf{c}(\mathbf{x}^k)_+\|_1/\xi_k \geq 1$ and there exists an infinite index set $\hat{\mathcal{K}} \subset \mathcal{K}$ such that $\mathcal{S}_\Psi^k = \emptyset$, for all $k \in \hat{\mathcal{K}}$;
- ii) $\liminf_{k \in \mathcal{K}} \|\mathbf{c}(\mathbf{x}^k)_+\|_1/\xi_k \geq 1$ and $\mathcal{S}_\Psi^k \neq \emptyset$, for all $k \in \mathcal{K}$ sufficiently large;
- iii) $\liminf_{k \in \mathcal{K}} \|\mathbf{c}(\mathbf{x}^k)_+\|_1/\xi_k < 1$.

Suppose i) holds. By Step 1, it follows that $\|\mathcal{P}(\mathbf{0} \mid \partial_{\epsilon_k} \|\mathbf{c}(\mathbf{x}^k)_+\|_1)\| \xrightarrow[k \in \hat{\mathcal{K}}]{} 0$. Consequently, since $\epsilon_k \rightarrow 0$, [40, Lemma 3.2 (iii)] guarantees that $\mathbf{0} \in \partial\|\mathbf{c}(\mathbf{x}^*)_+\|_1$.

Assume now that case ii) holds. This yields

$$\max_{k \in \mathcal{K}} \left\{ 1 - \frac{\|\mathbf{c}(\mathbf{x}^k)_+\|_1}{\xi_k}, 0 \right\} \xrightarrow[k \in \mathcal{K}]{} 0. \quad (21)$$

Additionally, by the way Algorithm 1 updates the value ρ_k , we must have $\rho_k > f(\mathbf{x}^*) + M$, for any large iteration k . Recalling that $\nu_k \downarrow 0$, $\epsilon_k \downarrow 0$ (because $\nu_{\text{opt}} = 0$ and $\epsilon_{\text{opt}} = 0$), $\mathcal{S}_{\Psi}^k \neq \emptyset$ for any large iteration $k \in \mathcal{K}$, and the result of Lemma 2.1, it follows, for any large $k \in \mathcal{K}$, that there exists $\{\mathbf{x}^{k,j}\}_{j=1}^{n+1} \subset \mathcal{B}(\mathbf{x}^k, \epsilon_k)$, and $\boldsymbol{\lambda}^k \in \mathbb{R}^{n+1}$ with $\mathbf{e}^T \boldsymbol{\lambda}^k = 1$, such that there exists $\mathbf{r}^k \in \partial \Psi_{\xi_k, \rho_k}(\mathbf{x}^k)$ with $\|\mathbf{r}^k\| \rightarrow 0$ and

$$\mathbf{r}^k \in \sum_{j=1}^{n+1} \lambda_j^k \left(\max \left\{ 1 - \frac{\|\mathbf{c}(\mathbf{x}^{k,j})_+\|_1}{\xi_k}, 0 \right\} \partial f(\mathbf{x}^{k,j}) + \sigma_j^k \partial \|\mathbf{c}(\mathbf{x}^{k,j})_+\|_1 \right), \quad (22)$$

where $\sigma_j^k \geq 1$, $j \in \{1, \dots, n+1\}$. Due to (21), $\|\mathbf{x}^k - \mathbf{x}^{k,j}\| \leq \epsilon_k$, and $\epsilon_k/\xi_k \rightarrow 0$, we get, for any $j \in \{1, \dots, n+1\}$,

$$\max \left\{ 1 - \frac{\|\mathbf{c}(\mathbf{x}^{k,j})_+\|_1}{\xi_k}, 0 \right\} \xrightarrow{k \in \mathcal{K}} 0.$$

The above limit together with (22) tell us that $\|\mathcal{P}(\mathbf{0} \mid \partial_{\epsilon_k} \|\mathbf{c}(\mathbf{x}^k)_+\|_1)\| \xrightarrow{k \in \mathcal{K}} 0$. So, since $\epsilon_k \downarrow 0$, we obtain $\mathbf{0} \in \partial \|\mathbf{c}(\mathbf{x}^*)_+\|_1$ due to [40, Lemma 3.2 (iii)].

Finally, we assume the validity of case iii), and we claim $\|\mathbf{c}(\mathbf{x}^k)_+\|_1 \rightarrow 0$. Indeed, because $\{\xi_k\}$ is a monotone decreasing and bounded sequence, it must follow that $\xi_k \downarrow \xi$, for some $\xi \geq 0$. According to Lemma 4.1, there exists $\{\delta_k\} \subset \mathbb{R}_+^*$ satisfying $\delta_k \rightarrow 0$ such that

$$\xi_k \leq \|\mathbf{c}(\mathbf{x}^k)_+\|_1 + \delta_k \Rightarrow 1 \leq \liminf_{k \in \mathcal{K}} \frac{\|\mathbf{c}(\mathbf{x}^k)_+\|_1}{\xi_k} + \liminf_{k \in \mathcal{K}} \frac{\delta_k}{\xi_k} \Rightarrow \liminf_{k \in \mathcal{K}} \frac{\delta_k}{\xi_k} > 0.$$

This yields $\bar{\xi} = 0$. However, by the way we have designed our algorithm, this only happens if there exists an infinite index set $\hat{\mathcal{K}} \subset \mathbb{N}$ with $\|\mathbf{c}(\mathbf{x}^k)_+\|_1 \xrightarrow{k \in \hat{\mathcal{K}}} 0$. Therefore, given any $\tau > 0$, there exists \hat{k}

such that $\|\mathbf{c}(\mathbf{x}^{\hat{k}})_+\|_1 < \tau$ and $\xi_{\hat{k}} < \tau$. Because of Lemma 4.3, this ensures $k \geq \hat{k} \Rightarrow \|\mathbf{c}(\mathbf{x}^k)_+\|_1 \leq \tau$. Since $\tau > 0$ is arbitrary, we obtain $\|\mathbf{c}(\mathbf{x}^k)_+\|_1 \rightarrow 0$, which guarantees that \mathbf{x}^* is a feasible point.

Then, Lemma 4.4 assures that, for all k large enough, we have $\mathcal{S}_{\Psi}^k \neq \emptyset$ and $\rho_k = \bar{\rho}$ for some $\bar{\rho} \in \mathbb{R}$. Notice that Step 5 ensures $f(\mathbf{x}^*) - \bar{\rho} \leq -M$. This information, together with Step 1, Lemma 2.2, the limit $\epsilon_k/\xi_k \rightarrow 0$ and the reasoning used in the proof of Theorem 2.2, prove that the cluster point \mathbf{x}^* satisfies the ϵ -ANOC. \square

Remark 3. *Some comments regarding Algorithm 1 are in order:*

- (a) *The technique based on sampling points to approximate $\mathcal{G}_\epsilon(\mathbf{x})$ that originated the method known as Gradient Sampling (GS) [16, 40], and its recent variants [24, 25, 33, 34, 35, 44] have shown to be effective tools for minimizing nonsmooth and nonconvex functions. The good results of those methods assure that Step 1 is not just an idealized step, but it can be performed in practice¹. In addition, BFGS techniques, which are not guaranteed to converge but work well in practice, may be used as well [22, 42]. Both possibilities were considered in Section 5.*
- (b) *Looking carefully at the proof of Theorem 4.1, we can understand the importance of Step 3 inside Algorithm 1. When $\|\mathbf{c}(\mathbf{x}^k)_+\|_1/\xi_k$ is greater than or approaching one, this implies that we are losing information about the objective function along the iterations, meaning that the method is tending to a stationary point of the infeasibility measure. Therefore, this suggests that one should choose $\omega \approx 1$.*
- (c) *For the case that problem (P) satisfies calmness [12] and ρ is large enough, every local minimizer of (P) is also a local minimizer of $\Psi_{\xi, \rho}$ for every ξ small enough (this can be easily proven by following the same reasoning used in the proof of [18, Proposition 6.4.3]). This ensures that, in many cases, PACNO will find a good approximation to the solution of the problem without needing to bring ξ_k down to zero.*

¹Assuming that f and c_i , $i \in \{1, \dots, p\}$, are continuously differentiable in full-measure open subsets of \mathbb{R}^n , all hypotheses required by the convergence theory of the GS method are satisfied.

5 Numerical results

This section has the goal to illustrate different properties of PACNO. Subsection 5.1 is devoted to show that, since our method is based on a legitimate necessary optimality condition (i.e., our convergence result does not depend on any kind of constraint qualification), PACNO may achieve good convergence to the solution of the nonsmooth optimization problem even in the absence of calmness. Additionally, we exhibit a practical example in which PACNO prevents the greediness phenomenon. Subsection 5.2 reveals that our method may converge even when the optimization problem does not satisfy all of our convergence hypotheses. Finally, Subsection 5.3 aims to illustrate that Step 1 of PACNO may accept different solvers to find a stationary point of $\Psi_{\xi,\rho}$. The tests were performed in a notebook DELL Latitude 7490, processor Intel Core i7-8650U, CPU 2.11GHz, with 16GB RAM (64-bit) using Matlab R2018a.

5.1 Nonsmooth optimization in the absence of calmness

In the previous sections, we have discussed the theoretical benefits of applying the penalization strategy used in $\Psi_{\xi,\rho}$ over the exact penalization approach to produce a practical necessary optimality condition. However, one can wonder if our penalization idea has any practical advantage when compared to the standard exact penalization procedure. Aiming to elucidate this matter, we present a simple nonsmooth optimization problem:

$$\min_{x_1, x_2} f(x_1, x_2) = \max\{x_1^3 - x_2, x_2\} \quad \text{s.t.} \quad c(x_1, x_2) = (x_1 - 10)^2 \leq 0. \quad (23)$$

The optimal solution of this problem is given by $\mathbf{x}^* = (10, 500)^T$ with its respective optimal value $f_* = 500$. Due to our choice of the constraint function that defines the feasible set $\mathcal{F} = \{(x_1, x_2) \in \mathbb{R}^2 : x_1 = 10\}$, the calmness property does not hold at the optimal point. This has a great impact on the behavior of methods that are based on the exact penalization approach, since this implies that one cannot solve the original problem with a finite penalty parameter [12]. In case this parameter value is allowed to go to infinity, this usually produces a sequence of points that has a poor precision on the optimal function value.

To illustrate how a method based on exact penalization behaves in the absence of calmness, we have solved the nonsmooth problem (23) using the SLQP-GS² (version 1.3) [23] in its default mode, with the exception that we have allowed the SLQP-GS algorithm to run 10^4 iterations in order to observe the optimal value precision that such a method is able to reach. In a similar way, the GS algorithm [16, 40] was chosen to be our internal solver for Step 1 – a version that we call PACNO_{GS} – and we have used the following parameter values for PACNO_{GS}: $\rho_1 = 0$, $M = 10$, $\theta_\xi = 0.5$, $\theta_\epsilon = 0.25$, $\theta_\nu = 0.5$, $\omega = 0.99$, $\xi_1 = \epsilon_1 = \nu_1 = 1$ and $\xi_{\text{opt}} = \epsilon_{\text{opt}} = \nu_{\text{opt}} = 10^{-8}$. In addition, for each call of GS, the sampling radii were set first as $\min\{10\epsilon_k, 1\}$ with optimality tolerance given by $\min\{10\nu_k, 1\}$, and then $\min\{\epsilon_k, 1\}$ as the final sampling radius with optimality tolerance $\min\{\nu_k, 1\}$. The results are shown in Figs. 4 and 5, so that one can follow the reached distributions of the infeasibility measure in (a), and of the relative errors ((b) in the domain, and (c) in the image space).

Because the SLQP-GS method relies on exact penalization and, additionally, its convergence theory is established only when the nonsmooth problem satisfies the calmness property at the solution, it was expected that the precision achieved by the method regarding the optimal function value would not be satisfactory. On the other hand, our algorithm possesses a convergence theory even in the absence of calmness, and, as anticipated, the precision obtained related to the optimality measure is considerably better. In addition, the PACNO_{GS} algorithm is able to keep the iterates closer to the feasible set.

Nevertheless, one can easily represent the feasible set \mathcal{F} in a manner that the calmness property is satisfied. Indeed, we can consider the constraint $|x_1 - 10| \leq 0$ instead of $(x_1 - 10)^2 \leq 0$ (see

²The code can be found in <http://coral.ise.lehigh.edu/frankecurtis/software/>

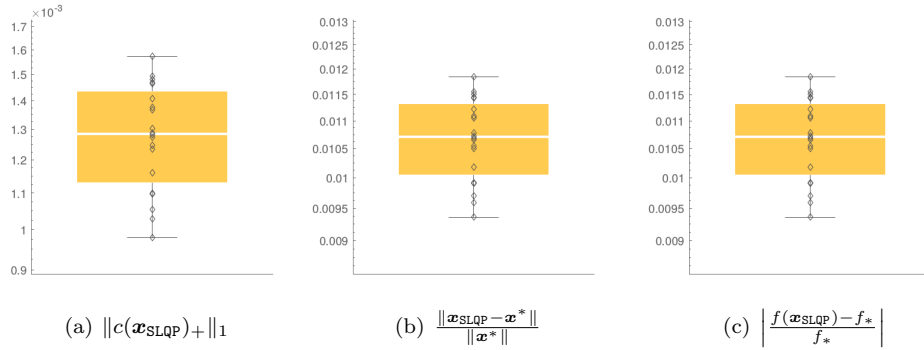


Figure 4: Boxplots of the results of twenty runs (depicted as \diamond) of the SLQP-GS method with different initial points chosen in a box $[-5, 5]^2$ centered at the optimal solution of problem (23). The last iterate obtained by the SLQP-GS method is represented by \mathbf{x}_{SLQP} .

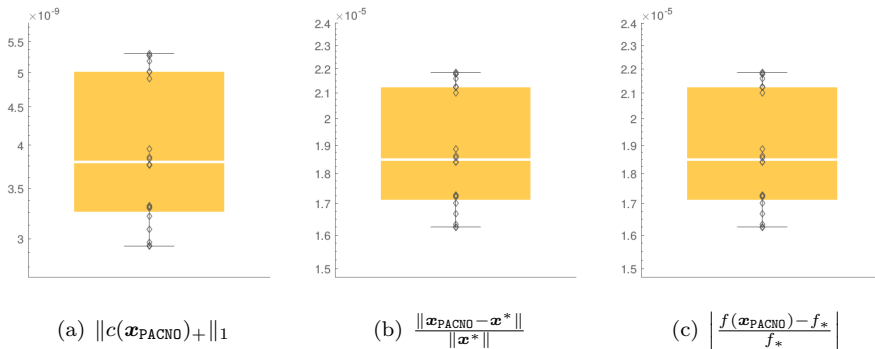


Figure 5: Boxplots of the result of twenty runs (depicted as \diamond) of the PACNO_{GS} method with different initial points chosen in a box $[-5, 5]^2$ centered at the optimal solution of problem (23). The last iterate obtained by the PACNO_{GS} method is represented by $\mathbf{x}_{\text{PACNO}}$.

Figs. 6 and 7). For this new optimization problem, the penalty parameter in the exact penalization approach no longer must go to infinity. However, methods based on exact penalization may experience another undesirable behavior: by initializing the method with an inappropriate value for the penalty parameter, such methods may present the greediness phenomenon. This anomaly can be seen when one tries to solve this new nonsmooth problem with the SLQP-GS method. In many runs, due to the initial value of the penalty parameter in the standard configuration of the SLQP-GS algorithm, the method gives excessive importance to the objective function in the first iterations, which carries the iterates away from the feasible set, preventing the method to converge even when one allows a large number of iterations to be performed (10^4 iterations). However, it is worth mentioning that, if the user sets a better scaled penalty parameter value, the SLQP-GS algorithm will easily converge to the solution point. In contrast, the PACNO_{GS} algorithm is able to reach the solution without needing to bring the parameter ξ down to zero (in the twenty runs that were performed, the mean value of ξ was kept above 10^{-3}), and, due to Lemma 4.3, the good behavior of PACNO_{GS} is not subjected to a tuned value of ξ .

5.2 Bilevel optimization via a nonsmooth approach

One of the most usual ways to solve a bilevel optimization problem [19] is to consider a mathematical problem with equilibrium constraints (MPEC) instead of the original multilevel instance. However,

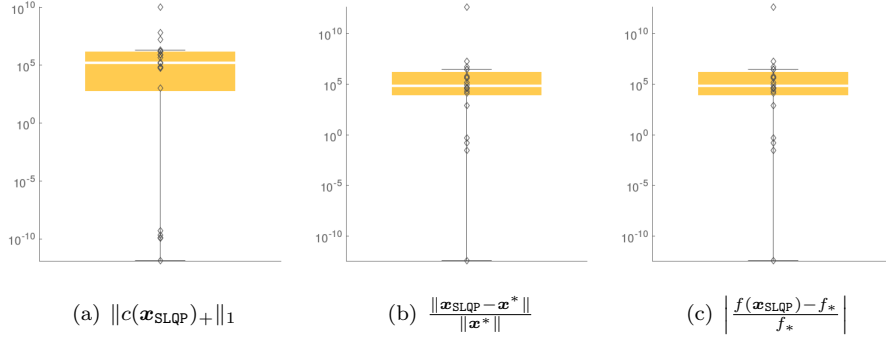


Figure 6: Boxplots of the result of twenty runs (depicted as \diamond) of the SLQP-GS method with different initial points chosen in a box $[-5, 5]^2$ centered at the optimal solution of problem (23). The last iterate obtained by the SLQP-GS method is represented by \mathbf{x}_{SLQP} .

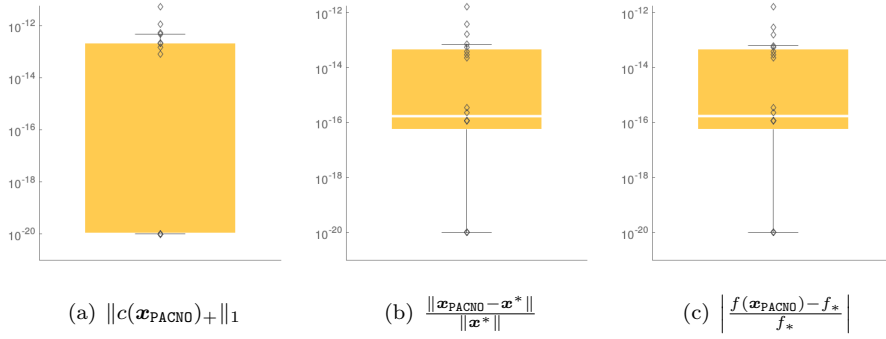


Figure 7: Boxplots of the result of twenty runs (depicted as \diamond) of the PACNO_{GS} method with different initial points chosen in a box $[-5, 5]^2$ centered at the optimal solution of problem (23). The last iterate obtained by the PACNO_{GS} method is represented by $\mathbf{x}_{\text{PACNO}}$.

there are situations in which this approach may generate an MPEC that is not equivalent to the original problem. For example, the authors of [26] show that the following minimization

$$\min_{x, \mathbf{y}} x \quad \text{s.t. } x \geq 0, \mathbf{y} \in \Lambda(x),$$

with $\Lambda(x) := \operatorname{argmin}_{\mathbf{y}} \{y_1^2 - y_2 \leq x, y_1^2 + y_2 \leq 0\}$, may generate an MPEC that does not possess an optimal solution when one considers the KKT conditions of the minimization problem related to the definition of $\Lambda(x)$. Indeed, consider the following nonsmooth single level optimization problem in which the first three constraints are associated with the KKT conditions of the implicit optimization problem presented in $\Lambda(x)$,

$$\begin{aligned} \min_{x, \mathbf{y}, \boldsymbol{\lambda}} \quad & x \\ \text{s.t.} \quad & |2\lambda_1 y_1 + 2\lambda_2 y_1 + 1| \leq 0 \\ & |-\lambda_1 + \lambda_2| \leq 0 \\ & \min\{-y_1^2 + y_2 + x, \lambda_1\}^2 + \min\{-y_1^2 - y_2, \lambda_2\}^2 \leq 0 \\ & -x \leq 0. \end{aligned} \tag{24}$$

One can see that, for every $x > 0$, it is always possible to find $\mathbf{y}(x)$ and $\boldsymbol{\lambda}(x)$ such that $(x, \mathbf{y}(x), \boldsymbol{\lambda}(x))$ is a feasible point. Although the objective function value of this problem converges to zero when $x \rightarrow 0$, the feasible set is empty for $x = 0$. This indicates that the MPEC instance (24) does not possess

an optimal point, and consequently, it cannot recover the original solution $\mathbf{p}^* := (x^*, y_1^*, y_2^*)^T = (0, 0, 0)^T$.

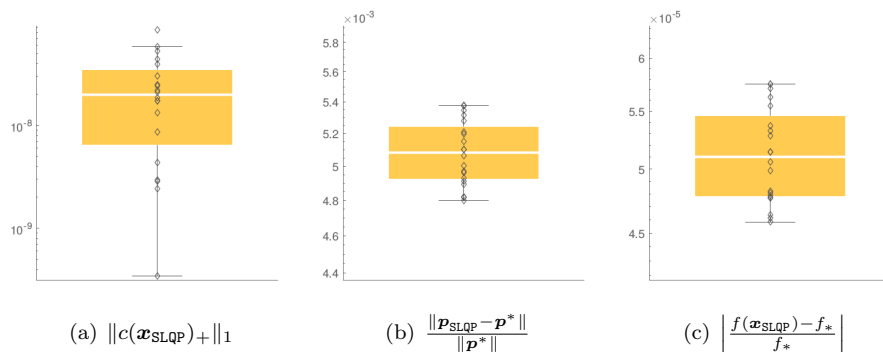


Figure 8: Boxplots of the results of twenty runs (depicted as \diamond) of the SLQP-GS method with different initial points chosen in a box $[-5, 5]^5$ centered at $(0, 0, 0, 50, 50)^T$. The primal variables (i.e., the first three coordinates) of the last iterate obtained by the SLQP-GS method is represented by \mathbf{p}_{SLQP} , whereas the complete vector is represented by \mathbf{x}_{SLQP} .

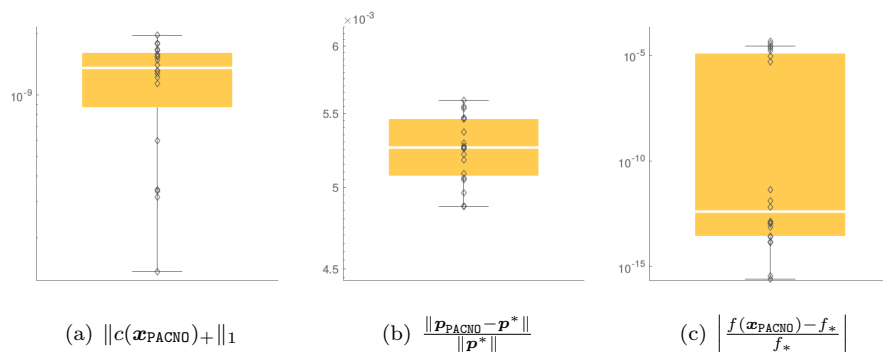


Figure 9: Boxplots of the results of twenty runs (depicted as \diamond) of the PACNO_{GS} method with different initial points chosen in a box $[-5, 5]^5$ centered at $(0, 0, 0, 50, 50)^T$. The primal variables (i.e., the first three coordinates) of the last iterate obtained by the PACNO_{GS} method is represented by $\mathbf{p}_{\text{PACNO}}$, whereas the complete vector is represented by $\mathbf{x}_{\text{PACNO}}$.

This type of problem does not fit into our theoretical assumptions, but there are reasons to believe that our approach may present a good behavior when applied to such a problem. Notice that, although the feasible set is empty when $x = 0$, slight perturbations on the constraints produce a nonsmooth optimization problem that accepts points for which the first coordinate is zero. Since the proposed sequential optimality conditions allow the use of first-order information at infeasible points (see (4)), one can expect that the PACNO_{GS} method (as defined in the previous subsection) might obtain a good approximation to the original solution of the bilevel optimization problem.

To verify our expectations, we have solved (24) using the PACNO_{GS} algorithm. Furthermore, we have also used the SLQP-GS algorithm as a way to have some comparative results. In the same way that it was done in the previous subsection, we allowed the SLQP-GS method to run 10^4 iterations in order to seek for the best precision that this method can achieve. The results are shown in Figs. 8 and 9.

Regarding the infeasibility measure, both methods are able to reach values close to zero – although, in the limit, the problem is infeasible. In contrast, when one looks to the optimality precision achieved by the algorithms, the PACNO_{GS} method presents a clear advantage in many runs. Because

the problem is infeasible, the **SLQP-GS** algorithm must bring the penalty parameter to a value that overshadows the objective function, giving too much importance to the infeasibility measure term. As a consequence, the method cannot substantially improve the optimality measure. On the other hand, since the **PACNO_{GS}** algorithm only occults the objective function when the iterates are too far from feasibility, the method is able to achieve high precision in the optimal value for many runs (13 out of 20).

5.3 The kissing number problem

In \mathbb{R}^n , how many non-overlapping spheres can touch, or kiss, simultaneously, another sphere of the same size? This quantity, known as *kissing number*, and here denoted by κ_n , is closely related to finding bounds for spherical codes and sphere packings [21]. Apparently dating back to 1694, when Isaac Newton and James Gregory disputed whether κ_3 was 12 or 13, finding kissing numbers are still open problems in most of the dimensions, for which just lower and upper bounds are known. Recent research has pursued improvements upon these bounds (e.g. [6, 43, 45]). For further details about the problem, including historical and mathematical related developments, we refer to the review [47] and the survey [10].

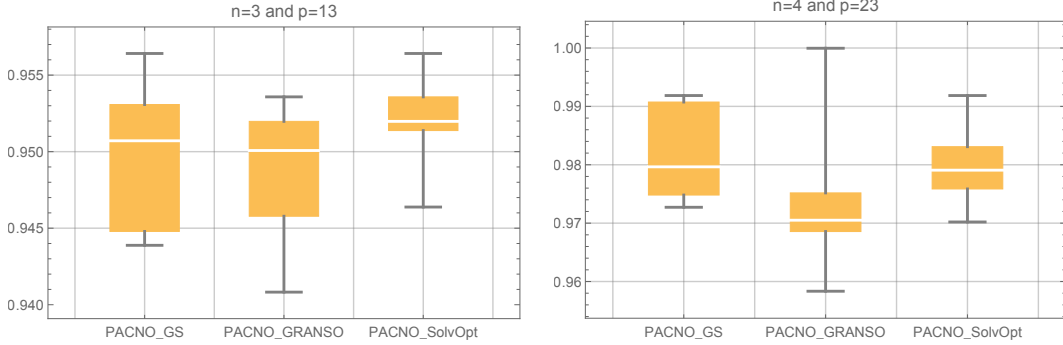


Figure 10: Boxplots of the maximum minimum distances obtained for (3, 13) and (4, 23).

Using known values, or bounds, of κ_n provide challenging instances for testing nonlinear programming algorithms [11, 41]. Indeed, for a given pair of integers (n, p) , and a given radius r , the following nonsmooth and nonconvex formulation aims at finding the centers $\mathbf{w}^i \in \mathbb{R}^n$, $i = 1, \dots, p$, of the spheres that touch a sphere centered at the origin

$$\max \min_{i \neq j} \|\mathbf{w}^i - \mathbf{w}^j\| \quad \text{s.t.} \quad \|\mathbf{w}^i\| = 2r, \quad i = 1, \dots, p.$$

Such a problem is equivalent to

$$\min \max_{i \neq j} (\mathbf{w}^i)^T \mathbf{w}^j \quad \text{s.t.} \quad \|\mathbf{w}^i\| = 2r, \quad i = 1, \dots, p, \quad (25)$$

having $n \times p$ variables, that may be arranged as $\text{vec}(\mathbf{w}^1 \cdots \mathbf{w}^p) \in \mathbb{R}^{np}$. To be addressed by Algorithm 1, each equality of (25) is turned into two inequalities, amounting to $2p$ inequality constraints. Setting the radius $r = 1/2$, for each choice of (n, p) , 20 initial points, were randomly and uniformly sampled in the box $[-2, 2]^{np}$. We have addressed nine instances, with (n, p) among the pairs: (3, 11), (3, 12), (3, 13), (4, 23), (4, 24), (4, 25), (5, 39), (5, 40), and (5, 41).

The corresponding instances were solved by Algorithm 1, using three possible solvers for computing the current approximation in Step 1, namely: **GS** [16, 40], **GRANSO** [22], and **SolvOpt** [37], respectively referred to as **PACNO_{GS}**, **PACNO_{GRANSO}**, and **PACNO_{SolvOpt}**.

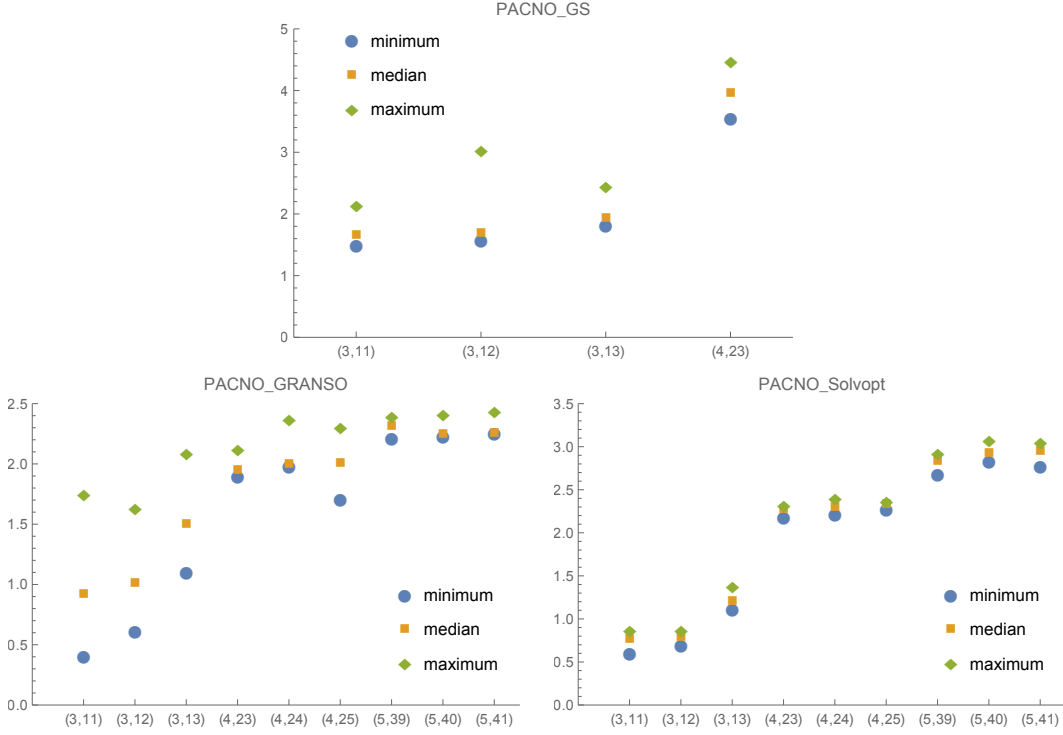


Figure 11: Range of the CPU time (in seconds, and \log_{10} scaled) demanded by PACNO for solving each instance, according with the inner solver: GS (top), GRANSO (bottom/left), and SolvOpt (bottom/right).

The algorithmic parameters were set as follows: $\rho_1 = 10$, $M = 5$, $\theta_\xi = 0.5$, $\theta_\epsilon = 0.25$, $\theta_\nu = 0.5$, $\omega = 0.99$, $\xi_1 = 10$, $\epsilon_1 = \nu_1 = 10^{-2}$, $\xi_{\text{opt}} = 10^{-8}$, $\epsilon_{\text{opt}} = 10^{-6}$, and $\nu_{\text{opt}} = 10^{-6}$. We have used tight tolerances aiming to reach accurate and better solutions. The maximum allowed number of iterations to be performed by the three inner solvers was set to 10^4 . Concerning their specific stopping criteria, we have made the following choices: for each call of GS, the sampling radii were set first as $\min\{10\epsilon_k, 1\}$ with optimality tolerance given by $\min\{10\nu_k, 1\}$, and then $\min\{\epsilon_k, 1\}$ as the final sampling radius with optimality tolerance $\min\{\nu_k, 1\}$; for GRANSO, the stationarity tolerance was $\tau_\diamond = \min\{\nu_k, 10^{-12}\}$, and the violation tolerance was $\tau_v = \min\{\epsilon_k, 10^{-6}\}$; and for SolvOpt, the tolerances for the relative error in the domain, and in the image space were set as $10^{-6}/\sqrt{n \cdot p}$, and $10^{-8}/\sqrt{n \cdot p}$, respectively.

The results for instances (3,11) and (3,12) were quite stable, without variability between the smallest and the largest values for the maximum minimum distances. In fact, the value 1.05146222 has been attained, no matter the initialization and the inner solver used. For (3,13), however, we have observed not only variability among the 20 initializations, but also slight differences among the reached values, as depicted in the distributions of Figure 10, and detailed in Table 1. It is worth mentioning that the values attained by PACNO_{SolvOpt} for (3,13) are even better than those reported in [11], corresponding to results obtained by differentiable solvers that addressed a smooth reformulation of (25).

As the time demanded by GS increases significantly with the problem dimension, instance (4,23) was the largest one addressed by the three inner solvers, and from (4,24) onwards, just GRANSO and SolvOpt were considered. The range of the CPU demanded by the inner solvers can be put in perspective by means of Figure 11, in which \log_{10} scaled values are presented. The plots evince not only the larger amount of time required by GS for solving instance (4,23), comparatively with the other two solvers, but also the wider variability reached by PACNO_{GRANSO} in comparison with

PACNO_{Solv0pt} for each instance.

Taking the values reported in [11] as reference, we have also noticed that, as the problem dimension increases, the quality of the obtained results slightly deteriorates (cf. Table 2). Moreover, the results reached by PACNO_{Solv0pt} were usually better than those obtained by PACNO_{GRANSO}, at the price of demanding more CPU time for being computed, as shown in Figure 11. The investigation of suitable parameter settings aiming to increase the attained maximum minimum distances will be subject of future research.

instance	algorithm	minimum	median	maximum	average
(3, 13)	PACNO _{GS}	0.9438814	0.9507097	0.9564136	0.9494221
	PACNO _{GRANSO}	0.9408229	0.9500742	0.9535789	0.9490796
	PACNO _{Solv0pt}	0.9463817	0.9519782	0.9564136	0.9522653
(4, 23)	PACNO _{GS}	0.9727178	0.9796446	0.9918580	0.9815232
	PACNO _{GRANSO}	0.9583368	0.9705020	0.9999590	0.9723568
	PACNO _{Solv0pt}	0.9702008	0.9790741	0.9918580	0.9797924

Table 1: Distribution values of the maximum minimum distances between two centers.

instance	algorithm	minimum	median	maximum	average
(4, 24)	PACNO _{GRANSO}	0.9434813	0.9578637	0.9730925	0.9579984
	PACNO _{Solv0pt}	0.9597049	0.9716582	0.9828751	0.9704148
(4, 25)	PACNO _{GRANSO}	0.9325337	0.9428726	0.9541111	0.9436839
	PACNO _{Solv0pt}	0.9474955	0.9554975	0.9617487	0.9552285
(5, 39)	PACNO _{GRANSO}	0.9431839	0.9561418	0.9647305	0.9550583
	PACNO _{Solv0pt}	0.9531385	0.9653934	0.9758015	0.9653913
(5, 40)	PACNO _{GRANSO}	0.9382251	0.9510471	0.9576960	0.9497185
	PACNO _{Solv0pt}	0.9451599	0.9615852	0.9727113	0.9604777
(5, 41)	PACNO _{GRANSO}	0.9320542	0.9453193	0.9547505	0.9439547
	PACNO _{Solv0pt}	0.9295390	0.9515950	0.9599363	0.9507821

Table 2: Distribution values of the maximum minimum distances between two centers (cont.)

6 Conclusion

We have proposed two sequential optimality conditions for a wide class of nonsmooth optimization problems. Both the weak ϵ -ANOC and ϵ -ANOC are legitimate necessary optimality conditions in the sense that they do not require any kind of constraint qualification to hold. In addition, when our stronger optimality condition is taken to the smooth context, we were able to prove that the ϵ -ANOC is stronger than AKKT condition and, moreover, that CAKKT and ϵ -ANOC are not connected to each other. Therefore, ϵ -ANOC is not just a new sequential optimality condition in the nonsmooth case, but it can also be seen as a novel necessary condition for smooth problems. Finally, we exhibited a practical algorithm able to generate both sequential optimality conditions as well as illustrative numerical results that highlight the potentialities of the devised algorithm.

Acknowledgements. *We would like to thank Michael L. Overton for hosting the third author during his internship at the New York University, Courant Institute, where much of this work was done, and also for his valuable comments. We are also in debt with Frank E. Curtis, for the algorithmic suggestions given in the early days of this study.*

This study was supported by Brazilian Funding Agencies Fundação de Amparo à Pesquisa do Estado de São Paulo - FAPESP (grants 2013/07375-0, 2013/05475-7, 2016/22989-2 and 2017/07265-0),

References

- [1] Andreani, R., Birgin, E.G., Martínez, J.M., Schuverdt, M.L.: Augmented Lagrangian methods under the constant positive linear dependence constraint qualification. *Mathematical Programming* **111**(1), 5–32 (2008)
- [2] Andreani, R., Haeser, G., Martinez, J.M.: On sequential optimality conditions for smooth constrained optimization. *Optimization* **60**(5), 627–641 (2011)
- [3] Andreani, R., Haeser, G., Secchin, L.D., Silva, P.J.S.: New sequential optimality conditions for mathematical problems with complementarity constraints and algorithmic consequences. *Optimization Online* pp. 1–30 (2018)
- [4] Andreani, R., Martnez, J., Svaiter, B.: A new sequential optimality condition for constrained optimization and algorithmic consequences. *SIAM Journal on Optimization* **20**(6), 3533–3554 (2010)
- [5] Andreani, R., Martnez, J.M., Schuverdt, M.L.: On the relation between constant positive linear dependence condition and quasinormality constraint qualification. *Journal of Optimization Theory and Applications* **125**(2), 473–483 (2005)
- [6] Bachoc, C., Vallentin, F.: New upper bounds for kissing numbers from semidefinite programming. *Journal of the American Mathematical Society* **21**(3), 909–924 (2008)
- [7] Birgin, E.G., Castelani, E.V., Martinez, A.L.M., Martínez, J.M.: Outer trust-region method for constrained optimization. *Journal of Optimization Theory and Applications* **150**(1), 142 (2011)
- [8] Birgin, E.G., Krejić, N., Martínez, J.M.: On the minimization of possibly discontinuous functions by means of pointwise approximations. *Optimization Letters* **11**(8), 1623–1637 (2017)
- [9] Bonnans, J.F., Gilbert, J.C., Lemaréchal, C., Sagastizábal, C.A.: *Numerical optimization: theoretical and practical aspects*, 2nd edn. Springer-Verlag Berlin Heidelberg (2006)
- [10] Boyvalenkov, P., Dodunekov, S., Musin, O.: A survey on the kissing numbers. *Serdica Mathematical Journal* **38**, 507–522 (2012)
- [11] Burachik, R., Kaya, Y.: An augmented penalty function method with penalty parameter updates for nonconvex optimization. *Nonlinear Analysis* **75**(3), 1158–1167 (2012)
- [12] Burke, J.: Calmness and exact penalization. *SIAM Journal on Control and Optimization* **29**(2), 493–497 (1991)
- [13] Burke, J.V.: An exact penalization viewpoint of constrained optimization. *SIAM Journal on Control and Optimization* **29**(4), 968–998 (1991)
- [14] Burke, J.V., Curtis, F.E., Lewis, A.S., Overton, M.L., Simões, L.E.A.: Gradient sampling methods for nonsmooth optimization. *ArXiv e-prints* (2018)
- [15] Burke, J.V., Lewis, A.S., Overton, M.L.: Approximating subdifferentials by random sampling of gradients. *Mathematics of Operations Research* **27**(3), 567–584 (2002)
- [16] Burke, J.V., Lewis, A.S., Overton, M.L.: A robust gradient sampling algorithm for nonsmooth, nonconvex optimization. *SIAM Journal on Optimization* **15**(3), 751–779 (2005)

- [17] Castelani, E.V., Martinez, A.L.M., Martínez, J.M., Svaiter, B.F.: Addressing the greediness phenomenon in nonlinear programming by means of proximal augmented Lagrangians. *Computational Optimization and Applications* **46**(2), 229–245 (2010)
- [18] Clarke, F.H.: *Optimization and nonsmooth analysis*, vol. 5. SIAM, Philadelphia, Canada (1990)
- [19] Colson, B., Marcotte, P., Savard, G.: An overview of bilevel optimization. *Annals of Operations Research* **153**(1), 235–256 (2007)
- [20] Conn, A.R., Gould, N.I.M., Toint, P.: A globally convergent augmented Lagrangian algorithm for optimization with general constraints and simple bounds. *SIAM Journal on Numerical Analysis* **28**(2), 545–572 (1991)
- [21] Conway, J.H., Sloane, N.J.C.: *Sphere Packings, Lattices and Groups*. Springer-Verlag, New York (1988)
- [22] Curtis, F.E., Mitchell, T., Overton, M.L.: A BFGS-SQP method for nonsmooth, nonconvex, constrained optimization and its evaluation using relative minimization profiles. *Optimization Methods and Software* **32**(1), 148–181 (2017)
- [23] Curtis, F.E., Overton, M.L.: A sequential quadratic programming algorithm for nonconvex, nonsmooth constrained optimization. *SIAM Journal on Optimization* **22**(2), 474–500 (2012)
- [24] Curtis, F.E., Que, X.: An adaptive gradient sampling algorithm for non-smooth optimization. *Optimization Methods and Software* **28**(6), 1302–1324 (2013)
- [25] Curtis, F.E., Que, X.: A quasi-Newton algorithm for nonconvex, nonsmooth optimization with global convergence guarantees. *Mathematical Programming Computation* **7**(4), 399–428 (2015)
- [26] Dempe, S., Dutta, J.: Is bilevel programming a special case of a mathematical program with complementarity constraints? *Mathematical Programming* **131**(1-2), 37–48 (2012)
- [27] Dutta, J., Deb, K., Tulshyan, R., Arora, R.: Approximate KKT points and a proximity measure for termination. *Journal of Global Optimization* **56**(4), 1463–1499 (2013)
- [28] Eremin, I.: The penalty method in convex programming. *Soviet Math. Dokl* **8**, 459–462 (1966)
- [29] Fasano, G., Liuzzi, G., Lucidi, S., Rinaldi, F.: A linesearch-based derivative-free approach for nonsmooth constrained optimization. *SIAM Journal on Optimization* **24**(3), 959–992 (2014)
- [30] Federer, H.: *Geometric measure theory*. Springer-Verlag, Berlin (1969)
- [31] Goldstein, A.A.: Optimization of Lipschitz continuous functions. *Mathematical Programming* **13**(1), 14–22 (1977)
- [32] Guo, L., Lin, G.H., Ye, J.J.: Second-order optimality conditions for mathematical programs with equilibrium constraints. *Journal of Optimization Theory and Applications* **158**(1), 33–64 (2013)
- [33] Helou, E.S., Santos, S.A., Simões, L.E.A.: On the differentiability check in gradient sampling methods. *Optimization Methods and Software* **31**(5), 983–1007 (2016)
- [34] Helou, E.S., Santos, S.A., Simões, L.E.A.: On the local convergence analysis of the gradient sampling method for finite max-functions. *Journal of Optimization Theory and Applications* **175**(1), 137–157 (2017)
- [35] Helou, E.S., Santos, S.A., Simões, L.E.A.: A fast gradient and function sampling method for finite-max functions. *Computational Optimization and Applications (Online)* (2018). URL <https://doi.org/10.1007/s10589-018-0030-2>
- [36] Hiriart-Urruty, J.B.: On optimality conditions in nondifferentiable programming. *Mathematical Programming* **14**(1), 73–86 (1978)

- [37] Kappel, F., Kuntsevich, A.V.: An implementation of Shors r-algorithm. *Computational Optimization and Applications* **15**(2), 193–205 (2000)
- [38] Kiwiel, K.C.: An exact penalty function algorithm for non-smooth convex constrained minimization problems. *IMA Journal of Numerical Analysis* **5**(1), 111–119 (1985)
- [39] Kiwiel, K.C.: Exact penalty functions in proximal bundle methods for constrained convex nondifferentiable minimization. *Mathematical Programming* **52**(1), 285–302 (1991)
- [40] Kiwiel, K.C.: Convergence of the gradient sampling algorithm for nonsmooth nonconvex optimization. *SIAM Journal on Optimization* **18**(2), 379–388 (2007)
- [41] Krejić, N., Martínez, J.M., Mello, M.P., Pilota, E.A.: Validation of an augmented Lagrangian algorithm with a Gauss-Newton Hessian approximation using a set of hard-spheres problems. *Computational Optimization and Applications* **16**(3), 247–263 (2000)
- [42] Lewis, A.S., Overton, M.L.: Nonsmooth optimization via quasi-Newton methods. *Mathematical Programming* **141**(1-2), 135–163 (2013)
- [43] Liberti, L.: Mathematical programming bounds for kissing numbers. In: A. Sforza, C. Sterle (eds.) *Optimization and Decision Science: Methodologies and Applications*, pp. 213–222. Springer (2017)
- [44] Loreto, M., Aponte, H., Cores, D., Raydan, M.: Nonsmooth Spectral Gradient Methods for Unconstrained Optimization. *EURO Journal on Computational Optimization* **5**(4), 529–553 (2017)
- [45] Machado, F.C., de Oliveira Filho, F.M.: Improving the semidefinite programming bound for the kissing number by exploiting polynomial symmetry. *Experimental Mathematics* pp. 1–8 (2017). URL <https://doi.org/10.1080/10586458.2017.1286273>
- [46] Nocedal, J., Wright, S.: *Numerical Optimization*, 2nd edn. Springer-Verlag New York (2006)
- [47] Pfender, F., Ziegler, G.M.: Kissing numbers, sphere packings, and some unexpected proofs. *Notices of the AMS* (September), 873–883 (2004)
- [48] Polak, E., Mayne, D., Wardi, Y.: On the extension of constrained optimization algorithms from differentiable to nondifferentiable problems. *SIAM Journal on Control and Optimization* **21**(2), 179–203 (1983)
- [49] Rockafellar, R.T., Wets, R., Wets, M.: *Variational analysis*, vol. 317. Springer (1998)
- [50] Zangwill, W.I.: Non-linear programming via penalty functions. *Management Science* **13**(5), 344–358 (1967)