

CONVERGENCE RATE FOR DIMINISHING STEPSIZE METHODS IN NONCONVEX CONSTRAINED OPTIMIZATION VIA GHOST PENALTIES

FRANCISCO FACCHINEI ^{a*}, VYACHESLAV KUNGURTSOV ^b,
LORENZO LAMPARIELLO ^c AND GESUALDO SCUTARI ^d

ABSTRACT. This is a companion paper to “Ghost penalties in nonconvex constrained optimization: Diminishing stepsizes and iteration complexity” (to appear in *Mathematics of Operations Research*). We consider the ghost penalty scheme for nonconvex, constrained optimization introduced in that paper, coupled with a diminishing stepsize procedure. Under an extended Mangasarian-Fromovitz-type constraint qualification we give an expression for the maximum number of iterations needed to achieve a given solution accuracy according to a natural stationarity measure, thus establishing the first result of this kind for a diminishing stepsize method for nonconvex, constrained optimization problems.

1. Introduction

We consider the nonconvex constrained optimization problem

$$\begin{aligned} & \underset{x \in \mathbb{R}^n}{\text{minimize}} && f(x) \\ & \text{s.t.} && g(x) \leq 0 \end{aligned} \tag{P}$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ are $C^{1,1}$ (i.e., continuously differentiable with locally Lipschitz gradients) functions on \mathbb{R}^n . Equality constraints can be added, but we avoid this for the sake of simplicity.

The main aim of this paper is to study the convergence rate of a Diminishing Stepsize Method (DSM) proposed by Facchinei *et al.* (2020). The analysis of Facchinei *et al.* (2020) fills a gap in the literature in that, for the first time, it shows convergence of a DSM for a constrained optimization problem with nonconvex constraints. The results in this paper complete and complement the analysis of Facchinei *et al.* (2020) by providing a convergence rate study for that method. We refer the interested reader to (Facchinei *et al.* 2020) for a more detailed discussion of all of the background and formulation of the method and its motivation.

This paper is dedicated with affection to Antonino Maugeri on the occasion of his 75th birthday

Our framework is of a generalized Sequential Quadratic Programming (SQP)-type. At each iteration x^v we generate a search direction $d(x^v)$ by solving a strongly convex optimization subproblem, which is described in Section 3, constructed along the lines discussed in the seminal papers of Burke (1989) and Burke and Han (1989) and also taking into account the developments of Scutari *et al.* (2014) and Facchinei *et al.* (2017). A step γ^v is then taken along $d(x^v)$ so that

$$x^{v+1} = x^v + \gamma^v d(x^v). \quad (1)$$

In (Facchinei *et al.* 2020), subsequential convergence to (generalized) stationary points of problem P is established for:

- (a) the classical diminishing stepsize method, in which the stepsize policy is chosen according to the following rules:

$$\lim_{v \rightarrow \infty} \gamma^v \downarrow 0 \quad \text{and} \quad \sum_{v=0}^{\infty} \gamma^v = \infty;$$

- (b) some more problem-oriented choices (including the case of a constant stepsize strategy) where the stepsize is held constant until some condition is met and then it is suitably reduced.

Iteration complexity for (b) is analyzed by Facchinei *et al.* (2020); in this paper we complete the analysis giving some corresponding complexity results for (a).

The paper is organized as follows. In the next section we recall some basic definitions. In Section 3 we define the search-direction subproblem and consider some related properties. In Section 4 we establish an upper bound on the number of iterations needed to satisfy a certain stopping criterion and also discuss the relationship of the stopping criterion to the KKT residual for problem (P).

2. KKT conditions and the eMFCQ

Consider Problem (P): our aim is to find a KKT point, i.e. a (feasible) point $x \in \mathbb{R}^n$ that, together with a vector of multipliers $\xi \in \mathbb{R}^m$, satisfies the KKT system

$$\begin{aligned} \nabla f(x) + \nabla g(x)\xi &= 0 \\ 0 \leq \xi \perp g(x) &\leq 0, \end{aligned} \quad (2)$$

where $\nabla f(x)$ is the gradient of f and $\nabla g(x)$ is the transposed Jacobian of g evaluated at x , and \perp mean orthogonal, so that $\xi \perp g(x)$ stands for $\xi^T g(x) = 0$. We find it convenient to describe the set of multipliers associated to a KKT point x as

$$M_1(x) \triangleq \left\{ \xi \mid \xi \in N_{\mathbb{R}_-^m}(g(x)), 0 = \nabla f(x) + \nabla g(x)\xi \right\},$$

where $N_{\mathbb{R}_-^m}(g(x))$ denotes the normal cone to the non positive orthant at $g(x)$ (note that this implies that $M_1(x)$ is surely empty at infeasible points, so that $M_1(x) \neq \emptyset$ implies feasibility for x). Note that, if x is feasible, i.e. $g(x) \leq 0$, condition $\xi \in N_{\mathbb{R}_-^m}(g(x))$ can be rewritten as

$$\xi_i \geq 0, \quad \xi_i g_i(x) = 0$$

for all i . Similarly we define the set of “abnormal” multipliers:

$$M_0(x) \triangleq \left\{ \xi \mid \xi \in N_{\mathbb{R}^m}(g(x) - \max_i \{g_i(x)_+\}e), 0 = \nabla g(x)\xi \right\}, \tag{3}$$

where $e \in \mathbb{R}^m$ is the vector with all components being one. Let \hat{x} be a local minimum point of (P), then it is well-known that either $M_1(\hat{x}) \neq \emptyset$, (the point is a KKT solution) or $M_0(\hat{x}) \neq \{0\}$ (the point is a Fritz-John point), or both. However, there may be feasible points x , possibly not minimum points, for which $M_0(x) \neq \{0\}$, these are precisely the points where the Mangasarian-Fromovitz constraint qualification fails. On the contrary, it is classical to show that if \hat{x} is not feasible, i.e. if $g_i(\hat{x}) > 0$ for at least an index $i \in \{1, \dots, m\}$, then the stationarity condition for the problem minimize $\max_x \{g_i(\hat{x})_+\}$, that is

$$0 \in \partial \max_i \{g_i(\hat{x})_+\},$$

is equivalent to $M_0(\hat{x}) \neq \{0\}$.

The Constraint Qualification (CQ) we use in this paper is the Mangasarian-Fromovitz CQ, suitably extended to infeasible points.

Definition 2.1. *We say that the extended Mangasarian-Fromovitz Constraint Qualification (eMFCQ) holds at x if*

$$M_0(x) = \{0\}.$$

Indeed, note that if x is feasible, then the condition $M_0(x) = \{0\}$ is nothing else but the standard MFCQ. We make the following blanket assumption

Assumption CQ The eMFCQ holds at any point $x \in \mathbb{R}^n$.

3. Direction Finding Subproblem and its Properties

At each iteration of our algorithm we move from the current iteration x^v along a direction $d(x^v)$ with a stepsize γ^v , see (1). While the stepsize is chosen according to classic diminishing stepsize rules, the direction $d(x^v)$ is the solution of a suitable strongly convex subproblem that we describe next.

Given a point x (which will actually be x^v in the algorithm) $d(x)$ is the unique solution of the optimization problem:

$$\begin{aligned} & \underset{d}{\text{minimize}} && \tilde{f}(d;x) \\ & \text{s.t.} && \tilde{g}(d;x) \leq \kappa(x)e \\ & && \|d\|_\infty \leq \beta, \end{aligned} \tag{P_x}$$

where $e \in \mathbb{R}^m$ is the vector with all components being one, $\kappa(x)$ a nonnegative quantity to be defined shortly, and β is a user-chosen positive constant. Moreover, \tilde{f} is a strongly convex surrogate of the original objective function f while \tilde{g} is a convex surrogate of the original constraint functions g (see Assumption A below for the conditions we impose on these surrogates).

The term $\kappa(x)e$ in the subproblem constraints serves to enlarge the feasible set of the subproblem in order to ensure it is always nonempty. The additional constraint $\|d\|_\infty \leq \beta$

allows one to avoid issues with search directions becoming too large. We denote by $\widetilde{\mathcal{X}}(x)$ the convex feasible set of subproblem (P_x) , i.e.

$$\widetilde{\mathcal{X}}(x) \triangleq \{d \in \mathbb{R}^n : \widetilde{g}(d;x) \leq \kappa(x)e, \|d\|_\infty \leq \beta\},$$

and, when convenient, we rewrite the constraint $\|d\|_\infty \leq \beta$ as $d \in \beta\mathbb{B}_\infty^n$, where \mathbb{B}_∞^n is the closed unit ball in \mathbb{R}^n associated with the infinity-norm.

The direction finding subproblem (P_x) is a direct generalization of the subproblems considered by Burke (1989), to which it reduces when the classical quadratic/linear approximations are used for \widetilde{f} :

$$\widetilde{f}(d;x) \triangleq \nabla f(x)^T d + \frac{1}{2} \|d\|^2; \quad \widetilde{g}(d;x) \triangleq g(x) + \nabla g(x)^T d. \quad (4)$$

Note that if these approximations are employed and we set $\kappa(x) = 0$ and $\beta = +\infty$, subproblem (P_x) boils down to the classical SQP-type subproblem. Here we adopt the approach of Burke (1989) by taking $\kappa(x)$ not necessarily zero and $\beta < +\infty$ in order to guarantee the existence and continuity of the solution mapping $d(x)$. In addition we introduce the use of general approximations \widetilde{f} and \widetilde{g} , this may be very convenient in practice by allowing flexibility in tailoring the direction finding subproblem to the problem at hand and to exploit any available specific structure in (P) , see Scutari *et al.* (2014) and Facchinei *et al.* (2017). Of course, an underlying assumption of our approach is that subproblem (P_x) can be solved efficiently. We do not insist on this point because it is very dependent on the choice of \widetilde{f} and \widetilde{g} which in turn is dictated by the original problem (P) . But we observe that the use of models that go beyond the standard quadratic/linear one in constrained optimization is emerging consistently in the literature, motivated, on the one hand, by the possibility to solve efficiently more complex subproblems than the classical quadratic ones and, on the other hand, by the desire of faster convergence rates, see for example the discussion in Section 3 of (Martínez 2017).

For this approach to be legitimate and lead to useful convergence results, we obviously need to make assumptions on the surrogate functions \widetilde{f} and \widetilde{g} .

Assumption Approx

Let O_d be an open neighborhood of $\beta\mathbb{B}_\infty^n$ and $\widetilde{f} : O_d \times \mathbb{R}^n \rightarrow \mathbb{R}$ and $\widetilde{g}_i : O_d \times \mathbb{R}^n \rightarrow \mathbb{R}$ for every $i = 1, \dots, m$ be continuously differentiable on O_d with respect to the first argument and such that

- A1):** $\widetilde{f}(\bullet;x)$ is a strongly convex function on O_d for every $x \in \mathbb{R}^n$ with modulus of strong convexity $c > 0$ independent of x ;
- A2):** $\widetilde{f}(\bullet;\bullet)$ is continuous on $O_d \times \mathbb{R}^n$;
- A3):** $\nabla_1 \widetilde{f}(\bullet;\bullet)$ is continuous on $O_d \times \mathbb{R}^n$;
- A4):** $\nabla_1 \widetilde{f}(0;x) = \nabla f(x)$ for every $x \in \mathbb{R}^n$;
- A5):** $\widetilde{g}_i(\bullet;x)$ is a convex function on O_d for every $x \in \mathbb{R}^n$;
- A6):** $\widetilde{g}_i(\bullet;\bullet)$ is locally Lipschitz continuous on $O_d \times \mathbb{R}^n$;
- A7):** $\widetilde{g}_i(0;x) = g_i(x)$ for every $x \in \mathbb{R}^n$;
- A8):** $\nabla_1 \widetilde{g}_i(\bullet;\bullet)$ is locally Lipschitz continuous on $O_d \times \mathbb{R}^n$;
- A9):** $\nabla_1 \widetilde{g}_i(0;x) = \nabla g_i(x)$, for every $x \in \mathbb{R}^n$;

where $\nabla_1 \tilde{f}(u; x)$ and $\nabla_1 \tilde{g}_i(u; x)$ denote the partial gradient of $\tilde{f}(\bullet; x)$ and $\tilde{g}_i(\bullet; x)$ evaluated at u .

These conditions are easily satisfied in practice and have been used in many recent papers; we refer the reader to (Scutari *et al.* 2014; Facchinei *et al.* 2017) as a good source of examples. Here we only note that the classical quadratic/linear approximations (4) satisfy Assumption A, provided f and g are $C^{1,1}$.

To complete the description of subproblem (P_x) , we must give the definition of $\kappa(x)$. Following Burke (1989), we set

$$\kappa(x) \triangleq (1 - \lambda) \max_i \{g_i(x)_+\} + \lambda \min_d \left\{ \max_i \{\tilde{g}_i(d; x)_+\} \mid \|d\|_\infty \leq \rho \right\}, \tag{5}$$

with $\lambda \in (0, 1)$ and $\rho \in (0, \beta)$. Note that the computation of $\kappa(x)$ requires the computation of the optimal value of the convex (see A5) optimization problem

$$\min_d \left\{ \max_i \{\tilde{g}_i(d; x)_+\} \mid \|d\|_\infty \leq \rho \right\}$$

that always has an optimal solution because the feasible set is nonempty and compact. This problem can easily be reformulated as a smooth optimization problem. In addition, if \tilde{g} is linear, as in the classical choice given in (4), this problem reduces to a linear program and can be efficiently solved. Note also that if x is feasible for (P), i.e. $g(x) \leq 0$, we have $\kappa(x) = 0$, so that our subproblem (P_x) is very similar to standard SQP subproblems. The term $\kappa(x)$ plays a key role when the SQP-type subproblems have an empty feasible set, a well-known issue with SQP schemes. In fact, observing that $\kappa(x)$ is always nonnegative being the sum of two nonnegative quantities, it restores feasibility by enlarging (with respect to the SQP choice $\tilde{g}(d; x) \leq 0$) the range of admissible values. Hence, the feasible set of problem (P_x) , for every x , is nonempty: choosing \hat{d} at which the minimum in the expression of $\kappa(x)$ is attained, we have

$$\tilde{g}(\hat{d}; x) \leq \min_d \left\{ \max_i \{\tilde{g}_i(d; x)_+\} \mid \|d\|_\infty \leq \rho \right\} e = \max_i \{g_i(\hat{d}; x)_+\} e,$$

and, in turn,

$$\begin{aligned} \tilde{g}(\hat{d}; x) &= (1 - \lambda) \tilde{g}(\hat{d}; x) + \lambda \tilde{g}(\hat{d}; x) \\ &\leq [(1 - \lambda) \max_i \{\tilde{g}_i(0; x)_+\} + \lambda \min_d \{\max_i \{\tilde{g}_i(d; x)_+\} \mid \|d\|_\infty \leq \rho\}] e = \kappa(x)e. \end{aligned}$$

Finally, we discuss briefly the KKT conditions for problem (P_x) . Observe preliminarily that the constraint $\|d\|_\infty \leq \beta$ corresponds to $2n$ bounds of the type $-\beta \leq d_i \leq \beta$. However, in what follows, we are interested only in the multipliers corresponding to the constraints $\tilde{g}(d; x) \leq 0$ and therefore we find it expedient to write the KKT conditions as

$$\begin{aligned} 0 &\in \nabla_1 \tilde{f}(d(x); x) + \nabla_1 \tilde{g}(d(x); x) \xi + N_{\beta \mathbb{B}_\infty^n}(d(x)) \\ 0 &\leq \xi \perp [\tilde{g}(d(x); x) - \kappa(x)] \leq 0 \end{aligned} \tag{6}$$

with the multipliers ξ satisfying the conditions $\xi \geq 0$ and $\xi^T \tilde{g}(d(x); x) = 0$, and where $N_{\beta \mathbb{B}_\infty^n}(d(x))$ is the normal cone to $\beta \mathbb{B}_\infty^n$ at $d(x)$.

The following proposition, collecting several useful facts about the objects considered so far, shows, among other things, that the KKT conditions are always satisfied at a solution of problem (P_x) .

Proposition 3.1. *Under Assumptions CQ and Approx, the following properties hold:*

- (a): $\kappa(\bullet)$ is locally Lipschitz continuous on \mathbb{R}^n ;
- (b): for every x , the feasible set $\widetilde{\mathcal{X}}(x)$ of (P_x) is non empty and (P_x) has a unique solution $d(x)$;
- (c): the MFCQ holds at every point of $\widetilde{\mathcal{X}}(x)$, for every x ;
- (d): the function $d(\bullet)$ is continuous on \mathbb{R}^n and $d(x) = 0$ if and only if x is a KKT solution for (P) ;
- (e): the unique solution $d(x)$ of (P_x) is a KKT solution of problem (P_x) and the set of KKT multipliers is locally bounded at any $x \in \mathbb{R}^n$.

Proof. Some of the statements are readily seen to hold and, in any case, with the exception of the fact that $d(x) = 0$ if and only if x is a KKT point of problem (P) , all the properties can immediately be derived from Propositions 2, 4, and 6 in (Facchinei *et al.* 2020), taking into account that here the eMFCQ is assumed to hold at every point. The statement about $d(x) = 0$ can be proved by simply comparing (2) and (6), taking into account A4 and A9 and the fact that $N_{\beta\mathbb{B}_\infty}(0) = \{0\}$: it is enough to observe that

$$\begin{aligned} \max_i \{g_i(x)_+\} - \kappa(x) &\leq \max_i \{g_i(x)_+\} - \max_i \{\tilde{g}_i(d(x); x)_+\} \\ &\leq \max_i \{g_i(x)_+\} - \max_i \{(g_i(x) + \nabla g_i(x)^T d(x))_+\} \quad (7) \\ &\leq \|\nabla g(x)^T\|_\infty \|d(x)\|, \end{aligned}$$

which is due to

$$\kappa(x) \geq \tilde{g}_i(d(x); x) \geq \tilde{g}_i(0; x) + \nabla \tilde{g}_i(0; x)^T d(x) = g_i(x) + \nabla g_i(x)^T d(x), \quad (8)$$

for every i . In turn, with $d(x) = 0$, $\max_i \{g_i(x)_+\} = \kappa(x)$ and the claim follows from Lemma 2 of (Facchinei *et al.* 2020). \square

4. Algorithm and Convergence Rate

Our scheme is based on the successive solution of subproblems (P_x) and the iterative process $x^{v+1} = x^v + \gamma^v d^v$, where the stepsizes γ^v are non-increasing and such that

$$\lim_{v \rightarrow \infty} \gamma^v \downarrow 0 \quad \text{and} \quad \sum_{v=0}^{\infty} \gamma^v = \infty. \quad (9)$$

Algorithm 1: DSM Algorithm for (P)

Data: $\gamma^0 \in (0, 1]$, $\delta \geq 0$, x^0 , $v \leftarrow 0$;

repeat

(S.1) compute $\kappa(x^v)$ and the solution $d(x^v)$ of problem (P_{x^v}) ;

(S.2) **if** $\|d^v\| \leq \delta$ **then**
 | **stop** and **return** x^v ;

end

(S.3) set $x^{v+1} = x^v + \gamma^v d(x^v)$, $v \leftarrow v + 1$;

end

The algorithm is always well defined if Assumption A, which guarantees existence and uniqueness of $d(x^v)$, holds. We also note that the stopping criterion at step S.2 is sound and sensible. In fact, in light of Proposition 3.1, $\|d(x)\|$ is a valid stationarity measure, being a continuous function which is zero if and only if x is a KKT point of problem (P). In subsection 4.1 we complete the study of this stationarity measure showing its quantitative relationship with the approximate KKT conditions residual.

Remark 4.1. *Note that we included the possibility of taking $\delta = 0$ in Algorithm 1. This choice will generally lead the algorithm to produce an infinite sequence $\{x^v\}$: in fact, the algorithm will stop only in the exceptional case of $d(x^v) = 0$ for some v , i.e. the algorithm will stop only if an exact KKT point is reached. This is the case considered by Facchinei et al. (2020): there, it is shown that, if $\delta = 0$, any limit point of the sequence $\{x^v\}$ is a KKT point for problem (P). Of course, to have finite termination, we must assume $\delta > 0$; however, the case $\delta = 0$ will be considered in order to state a boundedness condition used in Theorem 4.1, see Assumption Compact below. \square*

The distinctive aspect of Algorithm 1 is its simplicity, a feature shared with all DSMs. The main (and essentially only) computational burden is given by the computation of $\kappa(x^v)$ and of the solution for the strongly convex subproblem (P_{x^v}) . Before giving the result on the convergence rate properties of Algorithm 1, we state the last assumption needed.

Assumption Compact The sequence $\{x^v\}$ generated by Algorithm 1 with $\delta = 0$ is bounded. Hence, the whole sequence $\{x^v\}$ is contained in some compact set $S \subseteq \mathbb{R}^n$.

We will comment further on this assumption in Remark 4.3, after Theorem 4.1.

The proof of Theorem 4.1 below is based on the nonsmooth penalty function

$$W(x; \varepsilon) \triangleq f(x) + \frac{1}{\varepsilon} \max_i \{g_i(x)_+\} \tag{10}$$

where, as usual, $a_+ = \max\{0, a\}$. Essentially, and rather classically, we will show that for a certain sufficiently small value $\bar{\varepsilon}$ of the penalty parameter, $d(x^v)$ is a direction of sufficient decrease for $W(x; \bar{\varepsilon})$ at x^v . Based on this fact, we can derive the desired results. Note however that the penalty function itself is never used in the algorithm and, in particular, the user need not know the value of $\bar{\varepsilon}$, in contrast to standard penalty approaches in constrained optimization. For this reason, we use the term *ghost penalty*. We underline also the fact that we can see the penalty function as a Lyapunov function because we assume that the eMFCQ

holds; in (Facchinei *et al.* 2020), where this assumption is not made, the role played by the ghost penalty is more complicated.

Theorem 4.1. *Let $\{x^v\}$ be the sequence generated by Algorithm 1 with $\delta > 0$ under Assumptions CQ, Approx, and Compact. Then, Algorithm 1 finds a point for which $\|d(x^v)\| \leq \delta$ and stops in at most N iterations, where*

(a): *if γ^0 is sufficiently small, N is the first iteration index for which*

$$\sum_{v=0}^{N-1} \gamma^v \geq \frac{[W^0 - W^m]}{\omega \delta^2}, \tag{11}$$

where, $\omega = c/2$ (see Assumption Approx for the definition of c), $\bar{\epsilon}$ is a value of the penalty parameter to be specified in the proof, $W^0 \triangleq W(x^0; \bar{\epsilon})$, and W^m is the minimum value attained by the continuous function $W(x; \bar{\epsilon})$ on the compact set S ;

(b): *if no assumptions are made on γ^0 , N is the first iteration index for which*

$$\sum_{v=\bar{v}}^{\bar{v}+N-1} \gamma^v \geq \frac{[W^M - W^m]}{\omega \delta^2}, \tag{12}$$

for some suitable iteration \bar{v} and where W^M is the maximum value attained by the continuous function $W(x; \bar{\epsilon})$ on the compact set S .

Proof. Essentially, we aim at showing that a certain decrease in the penalty function is achieved at each iteration. The proof proceeds by a sequence of lemmas.

Lemma 4.1. *There is a positive constant Ξ such that for any $x \in S$ the multipliers ξ associated to the constraints $\tilde{g}(d; x) \leq \kappa(x)e$ at the solution $d(x)$ in subproblem (P_x) are bounded by Ξ , i.e. $\|\xi\|_\infty \leq \Xi$ for any optimal multiplier ξ .*

Proof. The lemma follows from Proposition 3.1 (e) and the compactness of S . □

The following lemma gives a bound on the directional derivative of the objective function f in the direction $d(x^v)$ at every iteration v .

Lemma 4.2. *At each iteration v the following bound holds for the directional derivative of the objective function at x^v in the direction $d(x^v)$:*

$$\nabla f(x^v)^T d(x^v) \leq -c\|d(x^v)\|^2 + m\Xi [\max_i \{g_i(x^v)_+\} - \kappa(x^v)]. \tag{13}$$

Proof. Recall first that, given x^v , $d(x^v)$ satisfies the KKT conditions (6) by Proposition 3.1:

$$0 \in \nabla_1 \tilde{f}(d(x^v); x^v) + \nabla_1 \tilde{g}(d(x^v); x^v) \xi^v + N_{\beta \mathbb{B}_{\infty}^m}(d(x^v)), \tag{14}$$

for some multipliers $\xi^v \in N_{\mathbb{R}_-^m}(\tilde{g}(d(x^v); x^v) - \kappa(x^v)e)$. Observe now that, thanks to A1 and A4,

$$\begin{aligned} \nabla_1 \tilde{f}(d(x^v); x^v)^T d(x^v) &= [\nabla_1 \tilde{f}(d(x^v); x^v) - \nabla_1 \tilde{f}(0; x^v) + \nabla_1 \tilde{f}(0; x^v)]^T d(x^v) \\ &\geq c\|d(x^v)\|^2 + \nabla f(x^v)^T d(x^v). \end{aligned} \tag{15}$$

Moreover, in view of A5,

$$-\nabla_1 \tilde{g}_i(d(x^v); x^v)^T d(x^v) \leq \tilde{g}_i(0; x^v) - \tilde{g}_i(d(x^v); x^v)$$

and, by A7, since ξ^v is nonnegative, in turn,

$$-\xi_i^v \nabla_1 \tilde{g}_i(d(x^v); x^v)^T d(x^v) \leq \xi_i^v [\tilde{g}_i(0; x^v) - \tilde{g}_i(d(x^v); x^v)] = \xi_i^v [g_i(x^v) - \kappa(x^v)], \quad (16)$$

where the equality follows observing that ξ^v belongs to $N_{\mathbb{R}^m}(\tilde{g}(d(x^v); x^v) - \kappa(x^v)e)$.

Therefore, by (14), (15) and (16), we have, for some $\zeta^v \in N_{\beta \mathbb{B}_\infty^m}(d(x^v))$,

$$\begin{aligned} c \|d(x^v)\|^2 + \nabla f(x^v)^T d(x^v) &\leq \nabla_1 \tilde{f}(d(x^v); x^v)^T d(x^v) \\ &= -\xi^{vT} \nabla_1 \tilde{g}(d(x^v); x^v)^T d(x^v) - \zeta^{vT} d(x^v) \\ &\leq \xi^{vT} [g(x^v) - \kappa(x^v)e] \leq \xi^{vT} [\max_i \{g_i(x^v)_+\} - \kappa(x^v)]e, \end{aligned}$$

and, thus, (13) easily follows taking into account Lemma 4.1 and the fact that, by definition, $\kappa(x^v) \leq \max_i \{g_i(x^v)_+\}$. □

Lemma 4.3. *For any positive ε , we have*

$$\begin{aligned} W(x^{v+1}; \varepsilon) - W(x^v; \varepsilon) &\leq \gamma^v \nabla f(x^v)^T d(x^v) - \frac{\gamma^v}{\varepsilon} \left[\max_i \{g_i(x^v)_+\} - \kappa(x^v) \right] \\ &\quad + \frac{(\gamma^v)^2}{2} \left(L_{\nabla f} + \frac{\max_i \{L_{\nabla g_i}\}}{\varepsilon} \right) \|d(x^v)\|^2, \end{aligned} \quad (17)$$

$$\begin{aligned} \nabla f(x^v)^T d(x^v) - \frac{1}{\varepsilon} \left[\max_i \{g_i(x^v)_+\} - \kappa(x^v) \right] &\leq -c \|d(x^v)\|^2 + \left(m\Xi - \frac{1}{\varepsilon} \right) \left[\max_i \{g_i(x^v)_+\} - \kappa(x^v) \right]. \end{aligned} \quad (18)$$

Proof. Inequality (17) is due to the following chain on relations:

$$\begin{aligned} W(x^{v+1}; \varepsilon) &- W(x^v; \varepsilon) \\ &= f(x^v + \gamma^v d(x^v)) - f(x^v) + \frac{1}{\varepsilon} \left[\max_i \{g_i(x^v + \gamma^v d(x^v))_+\} - \max_i \{g_i(x^v)_+\} \right] \\ &\stackrel{(a)}{\leq} \gamma^v \nabla f(x^v)^T d(x^v) + \frac{(\gamma^v)^2 L_{\nabla f}}{2} \|d(x^v)\|^2 \\ &\quad + \frac{1}{\varepsilon} \left[\max_i \{ (g_i(x^v) + \gamma^v \nabla g_i(x^v)^T d(x^v))_+ \} \right. \\ &\quad \left. - \max_i \{g_i(x^v)_+\} + \frac{(\gamma^v)^2 \max_i \{L_{\nabla g_i}\}}{2} \|d(x^v)\|^2 \right] \\ &\stackrel{(b)}{\leq} \gamma^v \nabla f(x^v)^T d(x^v) + \frac{1}{\varepsilon} \left[\max_i \{ (1 - \gamma^v) g_i(x^v)_+ + \gamma^v \kappa(x^v) \} - \max_i \{g_i(x^v)_+\} \right] \\ &\quad + \frac{(\gamma^v)^2}{2} \left(L_{\nabla f} + \frac{\max_i \{L_{\nabla g_i}\}}{\varepsilon} \right) \|d(x^v)\|^2 \\ &\leq \gamma^v \nabla f(x^v)^T d(x^v) - \frac{\gamma^v}{\varepsilon} \left[\max_i \{g_i(x^v)_+\} - \kappa(x^v) \right] \\ &\quad + \frac{(\gamma^v)^2}{2} \left(L_{\nabla f} + \frac{\max_i \{L_{\nabla g_i}\}}{\varepsilon} \right) \|d(x^v)\|^2, \end{aligned}$$

where (a) follows applying the descent lemma to f and g_i for every $i = 1, \dots, m$, with $L_{\nabla f}$ and $L_{\nabla g_i}$ being the Lipschitz moduli of ∇f and ∇g_i ; (b) holds for $\gamma^v \leq 1$ since, in view of (8), $\nabla g_i(x^v)^T d(x^v) \leq \kappa(x^v) - g_i(x^v)$. Inequality (18) is obtained by Lemma 4.2 by subtracting $\frac{1}{\bar{\epsilon}} \left[\max_i \{g_i(x^v)_+\} - \kappa(x^v) \right]$ from both sides of (13). \square

Proof of Theorem 4.1 (a). Set $\bar{\epsilon} \triangleq \frac{1}{m\Xi}$. We immediately get from (18)

$$\nabla f(x^v)^T d(x^v) - \frac{1}{\bar{\epsilon}} \left[\max_i \{g_i(x^v)_+\} - \kappa(x^v) \right] \leq -c \|d(x^v)\|^2.$$

Plugging this expression in (17) we obtain

$$\begin{aligned} W(x^{v+1}; \bar{\epsilon}) - W(x^v; \bar{\epsilon}) &\leq -\gamma^v c \|d(x^v)\|^2 + \frac{(\gamma^v)^2}{2} \left(L_{\nabla f} + \frac{\max_i \{L_{\nabla g_i}\}}{\bar{\epsilon}} \right) \|d(x^v)\|^2 \\ &= -\gamma^v \left[c - \frac{\gamma^v}{2} \left(L_{\nabla f} + \frac{\max_i \{L_{\nabla g_i}\}}{\bar{\epsilon}} \right) \right] \|d(x^v)\|^2. \end{aligned} \tag{19}$$

If

$$\gamma^v \leq \frac{c}{L_{\nabla f} + m\Xi \max_i \{L_{\nabla g_i}\}} \triangleq \bar{\gamma}, \tag{20}$$

the quantity in square brackets in the formula (19) is greater than or equal to $1/2$ and we can write

$$W(x^{v+1}; \bar{\epsilon}) - W(x^v; \bar{\epsilon}) \leq -\omega \gamma^v \|d(x^v)\|^2, \quad \omega \triangleq \frac{c}{2},$$

and hence

$$\gamma^v \|d(x^v)\|^2 \leq \frac{[W(x^v; \bar{\epsilon}) - W(x^{v+1}; \bar{\epsilon})]}{\omega}.$$

Assume now that $\gamma^0 \leq \bar{\gamma}$ and recall that the sequence $\{\gamma^v\}$ is non-increasing. Supposing that $\|d(x^v)\| > \delta$ for every $v \in \{0, \dots, N-1\}$ and taking the sum up to $N-1$, we have

$$\delta^2 \sum_{v=0}^{N-1} \gamma^v < \sum_{v=0}^{N-1} \gamma^v \|d(x^v)\|^2 \leq \frac{W(x^0; \bar{\epsilon}) - W(x^N; \bar{\epsilon})}{\omega} \leq \frac{W^0 - W^m}{\omega}.$$

We therefore see that the maximum number of iterations required in order to make $d(x^v)$ smaller than δ is N , with N such that

$$\sum_{v=0}^{N-1} \gamma^v \geq \frac{[W^0 - W^m]}{\omega \delta^2}$$

Proof of Theorem 4.1 (b). When considering a diminishing stepsize procedure, a finite \bar{v} exists such that, for every $v \geq \bar{v}$, (20) holds; clearly, the number of iterations \bar{v} is problem-dependent and relies on initial algorithmic choices such as the updating rule for the diminishing stepsize.

Supposing $\|d(x^v)\| > \delta$ for every $v \in \{0, \dots, \bar{v} + N\}$ and summing $\gamma^v \|d(x^v)\|^2 \leq [W(x^v; \bar{\epsilon}) - W(x^{v+1}; \bar{\epsilon})]/\omega$ from \bar{v} up to $\bar{v} + N - 1$, we have

$$\delta^2 \sum_{v=\bar{v}}^{\bar{v}+N} \gamma^v < \sum_{v=\bar{v}}^{\bar{v}+N} \gamma^v \|d(x^v)\|^2 \leq \frac{W(x^{\bar{v}}; \bar{\epsilon}) - W(x^{\bar{v}+N}; \bar{\epsilon})}{\omega} \leq \frac{W^M - W^m}{\omega},$$

Reasoning as in (a), we get (12). \square

Remark 4.2. In order for case (a) to hold, γ^0 must be “sufficiently small”, as stated in the theorem; the proof shows, more precisely, that actually γ^0 must be smaller than $\bar{\gamma}$, with $\bar{\gamma}$ defined in (20). \square

Remark 4.3. Meaningful results are obtained in Theorem 4.1 if the sequence $\{x^v\}$ generated by the algorithm with $\delta = 0$ is bounded. Although this is practically rather sensible, the question naturally arises if one can give a priori conditions guaranteeing the boundedness of the iterations. It is possible to give a satisfactory answer to these questions only for the price of a much more convoluted analysis; we eschew this for the sake of simplicity of presentation. The topic is however treated in great detail by Facchinei et al. (2020) to which we refer the reader for further information. All results on this aspect exposed by Facchinei et al. (2020) can be adapted to our setting. Here we only mention that there are two main avenues. One is that of giving explicit a priori coercivity-type conditions on problem (P) that guarantee boundedness of the iterates. The other approach is more direct and simply amounts to assuming that the original problem (P) includes a constraint of the type $x \in K$, where $K \subseteq \mathbb{R}^n$ is a (typically simple) compact, convex set. For example K could define upper and lower bounds on all variables. It is easy, although formally complicated, to show that all the results in this section still hold if we redefine subproblem (P_x), by appending the constraint K to it, as

$$\underset{d}{\text{minimize}} \tilde{f}(d; x^v) \quad \text{s.t.} \quad \tilde{g}(d; x^v) \leq \kappa(x^v) e, \quad \|d\|_\infty \leq \beta, \quad x^v + d \in K, \quad (21)$$

and by similarly modifying the feasible set of the minimization problem in the definition of κ . If one, then, requires the algorithm to start from an initial point $x^0 \in K$, this simple strategy obviously guarantees the boundedness of the iterations, which all belong to the compact set K , since $\gamma^v \in (0, 1]$ and K is convex. \square

The practical meaning of Theorem 4.1 depends on the actual rule according to which we choose γ^v . To make a concrete example, suppose that one relies on the classical generic term of the harmonic series and sets $\gamma^v = \frac{\gamma^0}{v+1}$ and that we are considering case (a) in Theorem 4.1. Observing that $\sum_{v=0}^{N-1} \gamma^v = \sum_{v=0}^{N-1} \frac{\gamma^0}{v+1} = \sum_{v=1}^N \frac{\gamma^0}{v}$ and $\gamma^0 \ln(N+1) < \sum_{v=1}^N \frac{\gamma^0}{v} < \gamma^0 [\ln N + 1]$,

$$\left| \exp \frac{|W^0 - W^m|}{\gamma^0 \omega \delta^2} - 1 \right| < N < \left\lceil \exp \frac{|W^0 - W^m|}{\gamma^0 \omega \delta^2} \right\rceil. \quad (22)$$

Case (b) is a variant of case (a). It is clear that if we use a diminishing stepsize rule, then, sooner or later, an iteration \bar{v} occurs for which $\gamma^{\bar{v}}$ satisfies the condition in point (a) and we can apply the results in (a) starting from that iteration. This case is worth considering because in practice it is the most realistic one, since in general it is difficult to establish the “sufficiently small” value $\bar{\gamma}$ that should be used in case (a). It is easy to see that when the generic term of the harmonic series $\frac{\gamma^0}{v+1}$ is employed, we can take

$$\bar{v} = \left\lceil \gamma^0 \frac{L_{\nabla f} + m \Xi \max_i \{L_{\nabla g_i}\}}{2c} \right\rceil - 1$$

and, since $\sum_{v=\bar{v}}^{\bar{v}+N-1} \gamma^v = \sum_{v=\bar{v}}^{\bar{v}+N-1} \frac{\gamma^0}{v+1} = \sum_{v=1}^N \frac{\gamma^0}{\bar{v}+v}$, we have

$$\gamma^0 \ln \left(\frac{\bar{v}+N}{\bar{v}+1} + \frac{1}{\bar{v}+1} \right) < \sum_{v=1}^N \frac{\gamma^0}{\bar{v}+v} < \gamma^0 \left[\ln \left(\frac{\bar{v}+N}{\bar{v}+1} \right) + \frac{1}{\bar{v}+1} \right],$$

and, in turn,

$$\left\lceil (\bar{v}+1) \exp \left[\frac{W^0 - W^m}{\gamma^0 \omega \delta^2} - \frac{1}{\bar{v}+1} \right] \right\rceil - \bar{v} < N < \left\lceil (\bar{v}+1) \exp \left[\frac{W^0 - W^m}{\gamma^0 \omega \delta^2} \right] \right\rceil - \bar{v}. \tag{23}$$

The discussion above hints at the fact that the use of a diminishing stepsize is, in a sense, a necessary evil, that one has to accept because the critical value $\bar{\gamma}$ is not known in advance. It is clear that, if $\gamma^0 \leq \bar{\gamma}$, it is better to have a sequence $\{\gamma^v\}$ that decreases “slowly”. From this point of view, the best possible choice would be the limiting case of a *constant stepsize* $\gamma^v = \gamma^0 \leq \bar{\gamma}$ for all v . Of course, this choice does not give rise to a DSM since the conditions (9) are not satisfied. Nevertheless, it is easy to check that this choice would still give termination in a finite number of steps. Indeed, since (20) is still satisfied, we can follow the proof of Theorem 4.1 from there and we still have $\|d(x^v)\|^2 \leq [W(x^v; \bar{\epsilon}) - W(x^{v+1}; \bar{\epsilon})] / \gamma^0 \omega$ for every v . Supposing that $\|d(x^v)\| > \delta$ for every $v \in \{0, \dots, N-1\}$ and taking the sum of iterations up to N , we have

$$N \gamma^0 \delta^2 < \sum_{v=0}^{N-1} \gamma^0 \|d(x^v)\|^2 \leq \frac{W(x^0; \bar{\epsilon}) - W(x^N; \bar{\epsilon})}{\omega} \leq \frac{W^0 - W^m}{\omega},$$

where $W^0 \triangleq W(x^0; \bar{\epsilon})$ and W^m is the minimum value attained by the continuous function $W(x; \bar{\epsilon})$ on the compact set S . Therefore, our procedure finds a point satisfying $\|d(x^v)\| \leq \delta$ in N iterations, with

$$N = \left\lceil \frac{[W^0 - W^m]}{\gamma^0 \omega \delta^2} \right\rceil. \tag{24}$$

This bound on the number of iterations is clearly much better in general than bounds like (22) or (23). Informally speaking, the bound (24) should be regarded as an ideal, limiting case for DSMs, a case obtained by taking the slowest possibly decreasing sequence past the necessary value $\bar{\gamma}$.

4.1. Approximate KKT Conditions. We already discussed how $\|d(x)\|$ is a valid stationarity measure that can be resorted to in order to devise a sensible stopping criterion. However, it is of interest to give quantitative relations between the stopping criterion $\|d(x)\| \leq \delta$ and the fulfillment of δ -approximate KKT conditions (see below for a precise definition). It turns out that the condition $\|d(x)\| \leq \delta$ entails the $\hat{\delta}$ -approximate KKT conditions, where $\hat{\delta}$ is equal to δ up to a fixed multiplicative factor. Lemma 4.4 below is the key tool in establishing the connection between the $\|d(x)\| \leq \delta$ stopping criterion and the approximate KKT conditions. In what follows, $L_{\nabla \tilde{f}}$, $L_{\nabla \tilde{g}}$ and $L_{\tilde{g}}$ denote the Lipschitz moduli of $\nabla_1 \tilde{f}(\bullet; \bullet)$, $\nabla_1 \tilde{g}(\bullet; \bullet)$ and $\tilde{g}(\bullet; \bullet)$ on $\beta \mathbb{B}_\infty^n \times S$.

Lemma 4.4. *Let $\{x^v\}$ be the sequence generated by Algorithm 1 with $\delta = 0$ under Assumptions CQ, Approx and Compact. Then, a common positive constant a exists such that, for*

every x^v , we have

$$\|\nabla f(x^v) + \nabla g(x^v)\xi^v\| \leq \left[L_{\nabla \tilde{f}} + \left(L_{\nabla \tilde{g}} + \frac{1}{\beta} \right) \Xi \right] \|d(x^v)\|, \quad (25)$$

$$\max_i \{g_i(x^v)_+\} \leq \left(\frac{L_{\tilde{g}}}{\lambda} + a \right) \|d(x^v)\|, \quad (26)$$

$$\max_i |g_i(x^v)\xi_i^v| \leq \max \left\{ 2L_{\tilde{g}}, \left(\frac{L_{\tilde{g}}}{\lambda} + a \right) \right\} \Xi \|d(x^v)\|, \quad (27)$$

for some $0 \leq \xi^v \in \mathbb{R}^m$.

Proof. As for the perturbed gradient of the Lagrangian-related condition, letting $\xi^v \in N_{\mathbb{R}^m}(\tilde{g}(d(x^v); x^v) - \kappa(x^v)e)$, which implies $\xi \geq 0$, we have

$$\begin{aligned} \|\nabla f(x^v) + \nabla g(x^v)\xi^v\| &= \|\nabla_1 \tilde{f}(0; x^v) - \nabla_1 \tilde{f}(d(x^v); x^v) + \nabla_1 \tilde{f}(d(x^v); x^v) \\ &\quad + \nabla_1 \tilde{g}(d(x^v); x^v)\xi^v + \nabla_1 \tilde{g}(0; x^v)\xi^v - \nabla_1 \tilde{g}(d(x^v); x^v)\xi^v\| \\ &\leq L_{\nabla \tilde{f}} \|d(x^v)\| + L_{\nabla \tilde{g}} \|\xi^v\| \|d(x^v)\| + \|\zeta^v\|, \end{aligned} \quad (28)$$

for some $\zeta^v \in N_{\beta \mathbb{R}_+^m}(d(x^v))$, where the equality is due to A4 and A9 and the inequality follows from (6). Since $\zeta^v = 0$ whenever $\|d(x^v)\|_\infty < \beta$, consider $d(x^v)$ such that $\|d(x^v)\|_\infty = \beta$: by (28),

$$\begin{aligned} \|\nabla f(x^v) + \nabla g(x^v)\xi^v\| &\leq L_{\nabla \tilde{f}} \|d(x^v)\| + L_{\nabla \tilde{g}} \|\xi^v\| \|d(x^v)\| + \|\zeta^v\| \frac{\|d(x^v)\|_\infty}{\beta} \\ &\leq L_{\nabla \tilde{f}} \|d(x^v)\| + L_{\nabla \tilde{g}} \|\xi^v\| \|d(x^v)\| + \frac{1}{\beta} \|\zeta^v\| \|d(x^v)\| \\ &\leq \left[L_{\nabla \tilde{f}} + \left(L_{\nabla \tilde{g}} + \frac{1}{\beta} \right) \Xi \right] \|d(x^v)\|, \end{aligned} \quad (29)$$

where the last inequality follows observing that, in view of Proposition 3.1 and Lemma 4.1 and by continuity, ζ^v is bounded on S , and we can take, without loss of generality, Ξ to be an upper bound also for $\|\zeta^v\|$. Regarding feasibility, by (7), we get $\max_i \{g_i(x^v)_+\} - \kappa(x^v) \leq L_{\tilde{g}} \|d(x^v)\|$, and, in turn,

$$\begin{aligned} \max_i \{g_i(x^v)_+\} - \kappa(x^v) &= \lambda [\max_i \{g_i(x^v)_+\} - \min_d \{\max_i \{\tilde{g}_i(d; x^v)_+\} \|d\|_\infty \leq \rho\}] \\ &\leq L_{\tilde{g}} \|d(x^v)\|. \end{aligned} \quad (30)$$

We now show that, for any x^v ,

$$\min_d \left\{ \max_i \{\tilde{g}_i(d; x^v)_+\} \|d\|_\infty \leq \rho \right\} \leq a \|d(x^v)\| \quad (31)$$

for some positive constant a .

Suppose on the contrary that subsequences $\{a^v\}_{\mathcal{N}} \in \mathbb{R}_+$ and $\{x^v\}_{\mathcal{N}}$ exist such that $a^v \xrightarrow{\mathcal{N}} +\infty$ and

$$\min_d \left\{ \max_i \{\tilde{g}_i(d; x^v)_+\} \|d\|_\infty \leq \rho \right\} > a^v \|d(x^v)\|. \quad (32)$$

Relation (32) implies

$$\min_d \left\{ \max_i \{ \tilde{g}_i(d; x^v)_+ \} \|d\|_\infty \leq \rho \right\} > 0. \tag{33}$$

Observing that, by A7, $\max_i \{g_i(x^v)_+\} \geq \min_d \{ \max_i \{ \tilde{g}_i(d; x^v)_+ \} \|d\|_\infty \leq \rho \} > 0$, thanks to Assumptions CQ and Compact, and in view of Lemma 2 of Facchinei *et al.* (2020), we have $\kappa(x^v) < \max_i \{g_i(x^v)_+\}$ and $\|d(x^v)\| \neq 0$ for every $v \in \mathcal{N}$. Moreover, since $\min_d \{ \max_i \{ \tilde{g}_i(d; \bullet)_+ \} \|d\|_\infty \leq \rho \}$ and $d(\bullet)$ are continuous on \mathbb{R}^n and S , respectively, see Proposition 3.1, we have, invoking Theorem 1 of Facchinei *et al.* (2020), $\|d(x^v)\| \xrightarrow{\mathcal{N}} \|d(\hat{x})\| = 0$, with \hat{x} cluster point of subsequence $\{x^v\}_{\mathcal{N}}$. Resorting again to Lemma 2 of Facchinei *et al.* (2020), $\kappa(\hat{x}) = \max_i \{g_i(\hat{x})\}_+$ and $\hat{d} \in \rho \mathbb{B}_\infty^n$ exist such that $\tilde{g}(\hat{d}; \hat{x}) < 0$. Hence, by continuity (A6), a neighborhood of \hat{d} exists such that, for any $d \in \rho \mathbb{B}_\infty^n$ belonging to it, and for $v \in \mathcal{N}$ sufficiently large, we have $\tilde{g}(d; x^v) < 0$, in contradiction to (33). Hence, (31) holds.

Combining (30) and (31), we get

$$\max_i \{g_i(x^v)_+\} \leq \frac{L_{\tilde{g}}}{\lambda} \|d(x^v)\| + \min_d \left\{ \max_i \{ \tilde{g}_i(d; x^v)_+ \} \|d\|_\infty \leq \rho \right\} \leq \left(\frac{L_{\tilde{g}}}{\lambda} + a \right) \|d(x^v)\|. \tag{34}$$

Concerning the complementarity condition, again without loss of generality, letting $\max_i |g_i(x^v) \xi_i^v| = |g_{\bar{i}}(x^v) \xi_{\bar{i}}^v| > 0$ for some $\bar{i} \in \{1, \dots, m\}$ such that $\tilde{g}_{\bar{i}}(d(x^v); x^v) = \kappa(x^v)$, we distinguish two cases.

If $g_{\bar{i}}(x^v) > 0$, thanks to the local boundedness of the set of KKT multipliers (recall Proposition 3.1 and Lemma 4.1), we have

$$|g_{\bar{i}}(x^v) \xi_{\bar{i}}^v| \leq \max_i \{g_i(x^v)_+\} |\xi_{\bar{i}}^v| \leq \left(\frac{L_{\tilde{g}}}{\lambda} + a \right) \Xi \|d(x^v)\|. \tag{35}$$

If, on the contrary, $g_{\bar{i}}(x^v) < 0$, in view of A7,

$$0 \leq \kappa(x^v) = \tilde{g}_{\bar{i}}(d(x^v); x^v) - \tilde{g}_{\bar{i}}(0; x^v) + \tilde{g}_{\bar{i}}(0; x^v) \leq L_{\tilde{g}} \|d(x^v)\| + g_{\bar{i}}(x^v) < L_{\tilde{g}} \|d(x^v)\|,$$

and, thus,

$$\begin{aligned} |g_{\bar{i}}(x^v) \xi_{\bar{i}}^v| &= |\tilde{g}_{\bar{i}}(0; x^v) - \tilde{g}_{\bar{i}}(d(x^v); x^v) + \tilde{g}_{\bar{i}}(d(x^v); x^v)| |\xi_{\bar{i}}^v| \\ &\leq (L_{\tilde{g}} \|d(x^v)\| + |g_{\bar{i}}(d(x^v); x^v)|) |\xi_{\bar{i}}^v| \\ &\leq 2L_{\tilde{g}} |\xi_{\bar{i}}^v| \|d(x^v)\| \leq 2L_{\tilde{g}} \Xi \|d(x^v)\|. \end{aligned} \tag{36}$$

□

We remark that the results in Lemma 4.4, which are obtained considering $\delta = 0$, are still valid for any x^v that is generated by Algorithm 1 with $\delta > 0$, since the sequence produced in the latter case results in nothing else but a truncation of the one generated whenever δ is set to zero.

We now define x^v to be a δ -approximate KKT point for problem (P) if each of the conditions in the KKT system is satisfied within a tolerance δ : more precisely,

$$\max \left\{ \|\nabla f(x^v) + \nabla g(x^v)^T \xi^v\|, \max_i |g_i(x^v)_+|, \max_i |g_i(x^v) \xi_i^v| \right\} \leq \delta.$$

Note that in principle we should also include the term $\|\xi^v\|$ in the max above, but this is not necessary here, because we already know that $\xi^v \geq 0$. It is interesting to highlight that this definition of approximate KKT point meets some important criteria for an approximate, or sequential, optimality condition (Andreani *et al.* 2010). The following theorem establishes the desired connection between the stopping criterion of Algorithm 1 and approximate KKT points; its proof follows easily from (25)-(27).

Theorem 4.2. *Let $\{x^v\}$ be the sequence generated by Algorithm 1 under Assumptions CQ, Approx and Compact. If, at iteration v of Algorithm 1, we have $\|d(x^v)\| \leq \delta$, then x^v is a δ' -approximate KKT point for problem P with*

$$\delta' \triangleq \delta \max \left\{ \left[L_{\nabla \bar{f}} + \left(L_{\nabla \tilde{g}} + \frac{1}{\beta} \right) \Xi \right], \left[\frac{L_{\tilde{g}}}{\lambda} + a \right], \left[\max \left\{ 2L_{\tilde{g}}, \left(\frac{L_{\tilde{g}}}{\lambda} + a \right) \right\} \Xi \right] \right\}.$$

Acknowledgments

Francisco Facchinei was partially supported by Progetto di Ateneo Distributed optimization algorithms for Big Data. Vyacheslav Kungurtsev was supported by the OP VVV project CZ.02.1.01/0.0/0.0/16_019/0000765 “Research Center for Informatics”. Lorenzo Lampariello was partially supported by the MIUR PRIN 2017 (grant 20177WC4KE). Gesualdo Scutari was partially supported by the NSF Grants CIF 1564044, CIF 1719205, and CMMI 1832688; and the ARO under Grant W911NF1810238.

References

- Andreani, R., Martínez, J. M., and Svaiter, B. F. (2010). “A new sequential optimality condition for constrained optimization and algorithmic consequences”. *SIAM Journal on Optimization* **20**(6), 3533–3554. DOI: [10.1137/090777189](https://doi.org/10.1137/090777189).
- Burke, J. V. (1989). “A sequential quadratic programming method for potentially infeasible mathematical programs”. *Journal of Mathematical Analysis and Applications* **139**(2), 319–351. DOI: [10.1016/0022-247X\(89\)90111-X](https://doi.org/10.1016/0022-247X(89)90111-X).
- Burke, J. V. and Han, S.-P. (1989). “A robust sequential quadratic programming method”. *Mathematical Programming* **43**(1), 277–303. DOI: [10.1007/BF01582294](https://doi.org/10.1007/BF01582294).
- Facchinei, F., Kungurtsev, V., Lampariello, L., and Scutari, G. (2020). “Ghost penalties in nonconvex constrained optimization: Diminishing stepsizes and iteration complexity”. *Mathematics of Operations Research*, to appear.
- Facchinei, F., Lampariello, L., and Scutari, G. (2017). “Feasible methods for nonconvex nonsmooth problems with applications in green communications”. *Mathematical Programming* **164**, 55–90. DOI: [10.1007/s10107-016-1072-9](https://doi.org/10.1007/s10107-016-1072-9).
- Martínez, J. M. (2017). “On high-order model regularization for constrained optimization”. *SIAM Journal on Optimization* **27**(4), 2447–2458. DOI: [10.1137/17M1115472](https://doi.org/10.1137/17M1115472).
- Scutari, G., Facchinei, F., Song, P., Palomar, D. P., and Pang, J.-S. (2014). “Decomposition by partial linearization: parallel optimization of multi-agent systems”. *IEEE Transactions on Signal Processing* **62**(3), 641–656. DOI: [10.1109/TSP.2013.2293126](https://doi.org/10.1109/TSP.2013.2293126).

-
- ^a University of Rome La Sapienza,
Department of Computer, Control, and System Engineering Antonio Ruberti,
Via Ariosto 25, 00185 Rome, Italy
- ^b Czech Technical University in Prague,
Department of Computer Science, Faculty of Electrical Engineering,
Prague, Czech Republic
- ^c Roma Tre University,
Department of Business Studies,
Via D'Amico 77, 00145 Rome, Italy
- ^d Purdue University,
School of Industrial Engineering,
West-Lafayette, IN, USA
- * To whom correspondence should be addressed | email: francisco.facchinei@uniroma1.it

Paper contributed to the meeting on "Variational Analysis, PDEs and Mathematical Economics", held in Messina, Italy (19–20 September 2019), on the occasion of Prof. Antonino Maugeri's 75th birthday, under the patronage of the *Accademia Peloritana dei Pericolanti*

Manuscript received 31 March 2020; published online 13 December 2020



© 2020 by the author(s); licensee *Accademia Peloritana dei Pericolanti* (Messina, Italy). This article is an open access article distributed under the terms and conditions of the [Creative Commons Attribution 4.0 International License](https://creativecommons.org/licenses/by/4.0/) (<https://creativecommons.org/licenses/by/4.0/>).