Improving the global convergence of Inexact Restoration methods for constrained optimization problems*

Roberto Andreani [†] Alberto Ramos[‡] Leonardo D. Secchin [§]
March 28, 2022 (revised June 26, 2024)

Abstract

Inexact restoration (IR) methods are an important family of numerical methods for solving constrained optimization problems with applications to electronic structures and bilevel programming among others areas. In these methods, the minimization is divided in two phases: decreasing infeasibility (feasibility phase) and improving optimality (optimality phase). The feasibility phase does not require the generated points to be feasible, so it has a practical appeal. In turn, the optimization phase involves minimize a simplified model of the problem over a linearization of the feasible set. In this paper, we introduce a new optimization phase through a novel linearization that carries more information about complementarity than the employed in previous IR strategies. We then prove that the resulting algorithmic scheme is able to converge globally to the so-called complementary approximate KKT (CAKKT) points. This global convergence result improves upon all previous results for this class of methods. In particular, convergence to KKT points is established with the very weak CAKKT-regularity condition. Furthermore, to the best of our knowledge, this is the first time that a method for general nonlinear programming has reached CAKKT points without exogenous assumptions. From the practical point of view, the new optimization phase does not require significant additional computational effort compared to the usual one. Our theory also provides new insights, even for the classical IR method, for cases where it is reasonable to compute exact feasible points in the feasibility phase. We present numerical experiments on CUTEst problems to support our findings.

Keywords: nonlinear optimization, inexact restoration, sequential optimality conditions, projection step

1 Introduction

We consider the optimization problem of the form

$$\min f(x) \quad \text{s.t.} \quad h(x) = 0, \quad g(x) \le 0 \tag{NLP}$$

^{*}This work has been partially supported by CEPID-CeMEAI (FAPESP 2013/07375-0), FAPESP (grants 2018/24293-0, 2017/18308-2 and 2023/08706-1), FAPES (grant 116/2019), CNPq (grants 306988/2021-6, 407147/2023-3 and 309136/2021-0), PRONEX - CNPq/FAPERJ (grant E-26/010.001247/2016) and ANID (Fondecyt grant 1231188).

[†]Department of Applied Mathematics, University of Campinas, Campinas, SP, Brazil. Email: andreani@ime.unicamp.br

[‡]Departamento de Matemática, Universidad de Tarapacá, casilla 7D, Arica, Chile. Email: aramosf@academicos.uta.cl

[§]Department of Applied Mathematics, Federal University of Espírito Santo, São Mateus, ES, Brazil. Email: leonardo.secchin@ufes.br

where $f: \mathbb{R}^n \to \mathbb{R}$, $h: \mathbb{R}^n \to \mathbb{R}^m$ and $g: \mathbb{R}^n \to \mathbb{R}^p$ are continuously differentiable functions. One class of numerical methods for solving (NLP) is the Inexact Restoration (IR). Roughly speaking, IR is an iterative method where its outer iteration is divided in two phases: feasibility phase and optimality phase. The first phase focuses on improving feasibility from the current point by applying a numerical method to a suitable infeasibility problem. In the second phase, a trial point with better optimality measure is computed, for example by minimizing a quadratic approximation of the objective function/Lagrangian over a linearization of the feasible set. Then, the generated trial point is accepted or rejected as a new iterate if some criterion is satisfied, commonly based on merit functions or filters, see [16, 20, 25, 27, 30]. IR methods have been successfully applied to problems where there is a specific way to improve feasibility. For instance, they were used for solving problems arising in electronic structure calculations [22], bilevel optimization [7], optimal control [28] and sample average approximation [12].

Classical global convergence results state that every feasible limit point of the sequence generated by some algorithm satisfies the so-called Karush-Kuhn-Tucker (KKT) conditions (Definition 2.1 below) if some constraint qualification (CQ) is valid at that point. The situation is more delicate when no CQ is valid. In this case, global convergence of several numerical methods has been established by means of sequential optimality conditions, which are satisfied at every local minimizer regardless of the validity of any CQ. See [8] and references therein. These conditions can be viewed as asymptotic versions of KKT and differ from each other essentially by the way the complementarity is described. Using them, we not only ensure that every feasible limit point generated by an algorithm satisfies a necessary optimality condition stronger than Fritz-John (FJ) [2], but we also guarantee the validity of KKT under very mild CQs; in particular, they do not imply the boundedness or uniqueness of the Lagrange multipliers. In this context, the global convergence of IR methods was established using the so-called approximate gradient projection (AGP) sequential optimality condition [31, 32].

In this work, we propose a new IR method for solving (NLP). Through a new optimization step, we improve the quality of the limit points generated by previous IR methods, showing that the new IR strategy converges to points associated with a strong sequential optimality condition. This new step is obtained by solving a quadratic subproblem that carries information about complementarity. Specifically, our focus is on the *complementary approximate* KKT (CAKKT) condition introduced in [6] (see Definition 2.2 below). CAKKT is considered one of the tightest sequential optimality condition as it implies others, including AGP [6, 8]. It provides the connection to the KKT conditions through a very mild CQ, namely CAKKT-regularity [5], that is strictly weaker than many others in the literature, such as *constant positive linear dependence* CQ (CPLD), *constant rank of the subspace component* CQ (CRSC) and *AKKT-regular* CQ (also known as *cone-continuity property* – CCP); see [5] for a complete overview. Furthermore, when (NLP) is a convex problem, CAKKT is necessary and sufficient to optimality [6, Theorem 4.2]. So, it is desirable to develop general-purpose algorithms that converge to CAKKT points while avoiding spurious non-local minimizers. This property makes such algorithms more preferable, at least theoretically, compared to other methods without this property; see Section 2.

Developing methods that converge to CAKKT points without exogenous assumptions is still a challenging task. For instance, it was shown in [6] that a safeguarded augmented Lagrangian (AL) method generates sequences whose feasible limit points fulfill the CAKKT condition if a certain measure of the infeasibility satisfies a generalized Kurdyka-Lojasiewicz (GKL) inequality. Under this assumption, the same holds for the primal-dual AL method considered in [8]. Some interior-point methods (IPM) are able to reach CAKKT points when only ordinary inequality constraints are present [24], but this capability is lost when inequality constraints are rewritten by inserting slack variables [8], as is commonly done in IPMs. The main feature of the proposed

IR method is that it reaches CAKKT points without imposing the GLK inequality or any exogenous assumptions, that is, its global convergence analysis is addressed only using standard hypotheses in the context of IR methods. To the best of our knowledge, this is the first method with this property. The key is to reformulate the CAKKT condition by means of gradient projections onto a suitable convex set that naturally fits within the IR framework. Furthermore, we introduce a sightly more flexible line-search procedure for accepting new iterates than what has traditionally been used.

This paper is organized as follows. Section 2 recalls some basic definitions and the main sequential optimality conditions for (NLP). We then present a reformulation of the CAKKT condition via gradient projections in Section 3. In Section 4, we present our IR algorithm in detail. Its theoretical convergence is addressed in Section 5. Computational experiments are presented in Section 6. Finally, Section 7 brings our conclusions and future research.

2 Preliminaries and notation

Our notation is standard in optimization and variational analysis. $\|\cdot\|$ denotes the Euclidean norm of a vector. We use \mathbb{R}^n_+ (respectively \mathbb{R}^n_-) to denote the subset of vectors in \mathbb{R}^n with non-negative (respectively non-positive) components. Given $a \in \mathbb{R}^m$, we set $a^- := \min\{a, 0\}$ and $a^+ := \max\{a, 0\}$ understood component-wise. Given a smooth function $q : \mathbb{R}^s \to \mathbb{R}^r$, $\nabla q(u)$ is $r \times s$ the matrix whose columns are $\nabla q_j(u)$, $j = 1, \ldots, s$. We denote the vector of ones by 1. The orthogonal projection of $u \in \mathbb{R}^n$ onto the closed convex set C is denoted by $\operatorname{proj}_C(u)$.

The Lagrangian function associated with (NLP) is

$$\mathcal{L}(x, \lambda, \mu) := f(x) + h(x)^T \lambda + g(x)^T \mu,$$

where $\lambda \in \mathbb{R}^m$ and $\mu \in \mathbb{R}^p_+$ are the dual variables. The feasible set of (NLP) is denoted by $\Omega = \{x \in \mathbb{R}^n \mid h(x) = 0, \ g(x) \leq 0\}$. The set of indices of active inequality constraints is denoted by $A(x) = \{j \in \{1, \dots, p\} \mid g_j(x) = 0\}$. For a given set-valued mapping $\mathcal{F} : \mathbb{R}^s \rightrightarrows \mathbb{R}^n$, the sequential Painlevé-Kuratowski outer/upper limit of $\mathcal{F}(u)$ as $u \to u^*$ [34] is defined as

$$\limsup_{u \to u^*} \mathcal{F}(u) = \{ y^* \in \mathbb{R}^n \mid \exists (u^k, y^k) \to (u^*, y^*) \text{ with } y^k \in \mathcal{F}(u^k), \ \forall k \in \mathbb{N} \}.$$

2.1 Review of sequential optimality conditions

We start our discussion with the definition of the KKT conditions.

Definition 2.1. We say that KKT conditions hold at the feasible point \bar{x} for (NLP) if there exist $\lambda \in \mathbb{R}^m$ and $\mu \in \mathbb{R}^p_+$ such that

$$\nabla f(\bar{x}) + \nabla h(\bar{x})\lambda + \nabla g(\bar{x})\mu = 0 \quad \text{and} \quad \min\{-g_j(\bar{x}), \mu_j\} = 0, \quad \forall j.$$
 (1)

In this case, we say that \bar{x} is a KKT point, and λ , μ are the multipliers.

Unlike the KKT conditions, sequential optimality conditions are genuine optimality conditions, which enables establishing global convergence of methods without mentioning CQs a priori. Also, they provide connections with KKT under very mild CQs. See for example [3, 5].

The complementarity in (1) can be written in different ways. The above form is related to the approximate KKT (AKKT) condition [2], which we recall next. We say that a feasible \bar{x}

satisfies the AKKT condition if there exist sequences $\{x^k\} \subset \mathbb{R}^n$, $\{\lambda^k\} \subset \mathbb{R}^m$ and $\{\mu^k\} \subset \mathbb{R}^p$ with x^k converging to \bar{x} such that

$$\nabla f(x^k) + \nabla h(x^k)\lambda^k + \nabla g(x^k)\mu^k \to 0 \quad \text{and} \quad \min\{-g_j(x^k), \mu_j^k\} \to 0, \quad \forall j.$$
 (2)

Through AKKT it was possible to improve the global convergence of several primal-dual methods, including AL, IPM and SQP algorithms, see [4] and references therein. In the context of IR methods, it was proposed in [32] the approximate gradient projection (AGP) condition, which holds at a feasible point \bar{x} if there exists a sequence $\{x^k\} \subset \mathbb{R}^n$ converging to \bar{x} such that

$$\operatorname{proj}_{\mathbb{L}_{\Omega}(x^k)}(-\nabla f(x^k)) \to 0,$$

where

$$\mathbb{L}_{\Omega}(x) := \left\{ d \in \mathbb{R}^n \,\middle|\, \begin{array}{c} \nabla h_i(x)^T d = 0, & \text{for all } i \\ g_i^-(x) + \nabla g_j(x)^T d \le 0, & \text{for all } j \end{array} \right\}. \tag{3}$$

Note that, in comparison with AKKT, AGP does not explicitly involve multipliers. Furthermore, AGP implies AKKT [2]. Regrettably, both conditions can lead us to accept spurious solution candidates as the next example shows, and therefore stronger conditions are desirable.

Example 2.1. Let us consider the example from [6]

$$\min \frac{(x_2 - 2)^2}{2} \quad s.t. \quad x_1 = 0, \quad x_1 x_2 = 0.$$

The unique minimizer is (0,2). In [6], it was stated that (0,1) is an AGP point. We affirm that every point $(0,\delta)$, $\delta \in \mathbb{R}$, is AGP (and thus, also AKKT). In fact, taking $\{x^k = (1/k,\delta)\}$ we have $\mathbb{L}_{\Omega}(x^k) = \{d \in \mathbb{R}^2 \mid d_1 = 0, \ \delta d_1 + d_2/k = 0\} = \{(0,0)\}$, and then $\operatorname{proj}_{\mathbb{L}_{\Omega}(x^k)}(-\nabla f(x^k)) = (0,0)$ for all k.

Next we recall the *complementary AKKT* (CAKKT) condition defined in [6].

Definition 2.2. We say that the CAKKT condition holds at the feasible point \bar{x} for (NLP) if there are sequences $\{x^k\} \subset \mathbb{R}^n$, $\{\lambda^k\} \subset \mathbb{R}^m$ and $\{\mu^k\} \subset \mathbb{R}^p$ such that (2) is valid, $\lambda_i^k h_i(x^k) \to 0, \forall i \text{ and } \mu_i^k g_i(x^k) \to 0, \forall j.$

Similar to KKT, when \bar{x} satisfies the AKKT/AGP/CAKKT condition we say that \bar{x} is an AKKT/AGP/CAKKT point, and that the corresponding sequence $\{x^k\}$ is an AKKT/AGP/CAKKT sequence. In [6], it was shown that any feasible limit point \bar{x} of a sequence generated by the safeguarded AL method is a CAKKT point provided that the measure of infeasibility $m(x) := \|h(x)\|^2 + \|g^+(x)\|^2$ associated with (NLP) satisfies the GKL inequality. This inequality ensures the existence of a continuous function ϕ that satisfies $\phi(\bar{x}) = 0$ and $\|m(x) - m(\bar{x})\| \le \phi(x)\|\nabla m(x)\|$ for every x near \bar{x} . Later on, a primal-dual AL method with good convergence properties was established [8], but it also requires further assumptions to ensure CAKKT points.

It is known that CAKKT implies AGP [6] (this is a direct consequence of the Theorem 2.1 below), and consequently AKKT. This indicates that an algorithm converging to CAKKT points is less likely to achieve non-minimizers than one that ensures only AGP. We revisit Example 2.1 to illustrate that this difference can be drastic.

Example 2.2 (Example 2.1 revisited). Let us consider Example 2.1, whose points $(0, \delta)$ are AGP. We affirm that only the global minimizer (0, 2) and (0, 0) are CAKKT points. In fact,

suppose that $\{x^k\}$ is a CAKKT sequence converging to $(0, \bar{x}_2)$, with associated dual sequence $\{(\lambda^k, \mu^k)\}$. Thus,

$$\begin{bmatrix} 0 \\ x_2^k - 2 \end{bmatrix} + \lambda^k \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \mu^k \begin{bmatrix} x_2^k \\ x_1^k \end{bmatrix} \to \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \qquad \lambda^k x_1^k \to 0, \qquad \mu^k x_1^k x_2^k \to 0.$$

Multiplying the second row of the first limit by x_2^k , we obtain $x_2^k(x_2^k-2) + \mu^k x_1^k x_2^k \to 0$. So the last limit gives $x_2^k(x_2^k-2) \to 0$, and thus $\bar{x}_2 = 2$ or $\bar{x}_2 = 0$. Since every feasible point is of the form $(0,\bar{x}_2)$, the preceding discussion implies that the only CAKKT points are (0,2) and (0,0).

Although AGP condition does not explicitly use any multipliers in its formulation, it is possible to reformulate it using them.

Theorem 2.1 ([8, Theorem 2.7]). Let \bar{x} be a feasible point of (NLP). Then, the AGP condition holds at \bar{x} if and only if there exist sequences $\{x^k\} \subset \mathbb{R}^n$, $\{\lambda^k\} \subset \mathbb{R}^m$ and $\{\mu^k\} \subset \mathbb{R}^p$ such that (2) holds and $\mu_j^k \min\{0, g_j(x^k)\} \to 0$.

The main difference between sequential optimality conditions lies in how the complementary condition is approximately satisfied. AKKT forces $\mu_j^k \to 0$ when $g_j(x^k) \not\to 0$ by requiring $\min\{\mu_j^k, -g_j(x^k)\} \to 0$ (see (2)). This condition allows the sequence $\{\mu_j^k\}$ to grow at any speed whenever $g_j(x^k) \to 0$. Besides the AKKT-like complementarity, the above theorem states that AGP also controls the behavior of μ^k by imposing $\mu_j^k \min\{0, g_j(x^k)\} \to 0$ for all j; so, the growth of $\{\mu_j^k\}$ is controlled by the way that $g_j(x^k)$ decreases from the infeasibility. CAKKT in turn imposes the most rigorous control on multipliers, that includes those on equality constraints, see Definition 2.2. In view of Theorem 2.1, a natural question is whether there is a sequential optimality condition based on the projected gradient that it is at least as strong as CAKKT. This question will be answered affirmatively in the next section.

3 CAKKT through projections

In this section we provide a reformulation of the CAKKT condition using projections, which is essential for the development of our new IR method. First, consider the non-empty convex set

$$\mathbb{L}_{\Omega}^{+}(x) := \left\{ (s, \ell, d) \in \mathbb{R}^{m} \times \mathbb{R} \times \mathbb{R}^{n} \left| \begin{array}{c} s_{i} h_{i}(x) + \nabla h_{i}(x)^{T} d = 0, \ \forall i \\ g_{j}^{-}(x) + \ell g_{j}^{+}(x) + \nabla g_{j}(x)^{T} d \leq 0, \ \forall j \end{array} \right\}.$$
 (4)

Note that $d \in \mathbb{L}_{\Omega}(x)$ implies $(0,0,d) \in \mathbb{L}_{\Omega}^+(x)$, but the "d-part" of these sets are generally different. For instance, consider the problem of minimizing $f(x_1,x_2) = (x_2-2)^2/2$ subject to $x_1 = 0, x_1x_2 \leq 0$. Clearly, $\mathbb{L}_{\Omega}(1/k,1) = \{0\} \times (-\infty,0]$ while

$$\mathbb{L}_{\Omega}^{+}(1/k,1) = \{(s,\ell,d) \in \mathbb{R} \times \mathbb{R} \times \mathbb{R}^2 \mid d_1 = -s/k, \ d_2 \le s - \ell\}.$$

It is straightforward to see that the sequence of projections of $(0,0,-\nabla f(1/k,1))=(0,0,0,1)$ onto this set converges to the non-null vector (1/3,-1/3,0,2/3). This is consistent with the fact that (0,1) is not a CAKKT point (Example 2.2). The next result shows that this is not a mere coincidence: in fact, CAKKT can be reinterpreted in terms of projections onto $\mathbb{L}_{\Omega}^+(x)$.

Theorem 3.1. Let \bar{x} be a feasible point of (NLP). Then, CAKKT holds at \bar{x} if and only if there is a sequence x^k converging to \bar{x} such that

$$\operatorname{proj}_{\mathbb{L}_{0}^{+}(x^{k})}(0,0,-\nabla f(x^{k})) \to (0,0,0). \tag{5}$$

Proof. Let us show that (5) implies CAKKT. By (5), there exists a sequence $x^k \to \bar{x}$ such that $(s^k, \ell_k, d^k) := \operatorname{proj}_{\mathbb{L}^+_{\Omega}(x^k)}(0, 0, -\nabla f(x^k)) \to (0, 0, 0)$. From the definition of orthogonal projection, (s^k, ℓ_k, d^k) is a global minimizer of

$$\min \frac{1}{2} \| -\nabla f(x^k) - d \|^2 + \frac{1}{2} \|s\|^2 + \frac{1}{2} \ell^2$$

$$s_i h_i(x^k) + \nabla h_i(x^k)^T d = 0, \ \forall i,$$
s.t.
$$g_j^-(x^k) + \ell g_j^+(x^k) + \nabla g_j(x^k)^T d \le 0, \ \forall j.$$

Since all the constraints are affine, the KKT conditions hold. Thus, there exist $\lambda^k \in \mathbb{R}^m$ and $\mu^k \in \mathbb{R}^p_+$ such that

$$\nabla f(x^k) + d^k + \nabla h(x^k)\lambda^k + \nabla g(x^k)\mu^k = 0, \tag{7a}$$

$$s_i^k + \lambda_i^k h_i(x^k) = 0, \quad \forall i, \qquad \ell_k + \sum_{j=1}^p \mu_j^k g_j^+(x^k) = 0,$$
 (7b)

$$\mu_j^k \left(g_j^-(x^k) + \ell_k g_j^+(x^k) + \nabla g_j(x^k)^T d^k \right) = 0, \quad \forall j,$$
 (7c)

for each k. Considering (7a) and $d^k \to 0$, we have that (2) holds. It follows from (7b) and $(s^k, \ell_k) \to (0, 0)$ that

$$\lambda_i^k h_i(x^k) \to 0 \text{ for all } i \text{ and } 0 \le \sum_{j=1}^p \mu_j^k g_j^+(x^k) \to 0.$$
 (8)

Due to the non-negativity of the terms in the above sum, the last limit implies $\mu_j^k g_j^+(x^k) \to 0$ for all j. Therefore, it is enough to show that $\mu_j^k g_j^-(x^k) \to 0$ for all j to ensure that $\mu_j^k g_j^-(x^k) \to 0$ for all j, and so conclude that \bar{x} is CAKKT. Doing the inner product of (7a) with d^k and using (7c), we arrive at (we eventually omit " (x^k) " when it is clear from the context)

$$\sum_{j=1}^{p} \mu_{j}^{k} g_{j}^{-} = -\ell_{k} \sum_{j=1}^{p} \mu_{j}^{k} g_{j}^{+} + \nabla f^{T} d^{k} + \sum_{i=1}^{m} \lambda_{i}^{k} \nabla h_{i}^{T} d^{k} + \|d^{k}\|^{2}$$

$$(9)$$

By the feasibility of (s^k, ℓ_k, d^k) (see (6)) we have that $\lambda_i^k \nabla h_i(x^k)^T d^k = -s_i^k \lambda_i^k h_i(x^k)$. Now, using the fact that $(s^k, \ell_k, d^k) \to (0, 0, 0)$ and (8) we have that all terms on the right side of (9) tend to zero. Since each term in the sum on the left-hand side of (9) is non-positive, they all also converge to zero. Therefore, $|\mu_j^k g_j(x^k)| = \max\{\mu_j^k g_j^+(x^k), -\mu_j^k g_j^-(x^k)\} \to 0$ for all j.

Conversely, suppose that \bar{x} is a CAKKT point. There exist sequences $x^k \to \bar{x}$, $\hat{\lambda}^k \in \mathbb{R}^m$ and $\hat{\mu}^k \in \mathbb{R}^p_+$ such that

$$w^{k} := \nabla f + \nabla h \widehat{\lambda}^{k} + \nabla g \widehat{\mu}^{k} \to 0, \quad \widehat{\lambda}_{i}^{k} h_{i} \to 0, \forall i, \quad \widehat{\mu}_{i}^{k} g_{i} \to 0, \forall j.$$
 (10)

For each k, set $(s^k, \ell_k, d^k) := \operatorname{proj}_{\mathbb{L}^+_{\Omega}(x^k)}(0, 0, -\nabla f(x^k))$. Since (s^k, ℓ_k, d^k) is an optimal solution of the projection problem, there exist $\lambda^k \in \mathbb{R}^m$, $\mu^k \in \mathbb{R}^p_+$ so that (7a)–(7c) are valid. Multiplying (7a) by d^k , the feasibility of (s^k, ℓ_k, d^k) , the inequality $\mu^k_j g^-_j(x^k) \leq 0$, (7b) and (7c) give

$$||d^k||^2 = -\nabla f^T d^k + \sum_{i=1}^m s_i^k \lambda_i^k h_i + \ell_k \sum_{j=1}^p \mu_j^k g_j^+ + \sum_{j=1}^p \mu_j^k g_j^- \le -\nabla f^T d^k - ||s^k||^2 - \ell_k^2,$$

which implies

$$||d^k||^2 + ||s^k||^2 + \ell_k^2 \le -\nabla f(x^k)^T d^k, \tag{11}$$

and so $-\nabla f(x^k)^T d^k \ge 0$. Now, we proceed to find an upper bound for $-\nabla f(x^k)^T d^k$. Note that $0 \le \|\nabla f(x^k) + d^k\|^2 = \|\nabla f(x^k)\|^2 + \|d^k\|^2 + 2\nabla f(x^k)^T d^k$ and (11) imply

$$2\|d^k\|^2 \le -2\nabla f(x^k)^T d^k \le \|\nabla f(x^k)\|^2 + \|d^k\|^2,$$

and thus $||d^k||^2 \leq ||\nabla f(x^k)||^2$. Therefore,

$$\|d^k\|^2 + \|s^k\|^2 + \ell_k^2 \leq \|\nabla f\|^2 + \|s^k\|^2 + \ell_k^2 \leq \|\nabla f\|^2 + \| - \nabla f - d^k\|^2 + \|s^k\|^2 + \ell_k^2 \leq 2\|\nabla f\|^2,$$

where in the last inequality we use the optimality of (s^k, ℓ_k, d^k) for the projection problem (6) and $(0,0,0) \in \mathbb{L}^+_{\Omega}(x^k)$. Thus, for every k, $\max\{\|d^k\|, \|s^k\|, |\ell_k|\} \leq 2\|\nabla f(x^k)\|$. Define $\zeta^k := g^-(x^k) + \ell_k g^+(x^k) + \nabla g(x^k)^T d^k$, $\forall k$. The feasibility of (s^k, ℓ_k, d^k) implies $\zeta^k \in \mathbb{R}^p_-$, $\forall k$, and hence $\zeta^k_j \widehat{\mu}^k_j \leq 0$, $\forall j, k$. Using (7a), (10) and the Cauchy-Schwarz inequality, we obtain

$$\begin{split} -\nabla f^T d^k &= -(w^k)^T d^k + \sum_{i=1}^m \widehat{\lambda}_i^k \nabla h_i^T d^k + \sum_{j=1}^p \widehat{\mu}_j^k \nabla g_j^T d^k \\ &= -(w^k)^T d^k - \sum_{i=1}^m \widehat{\lambda}_i^k s_i^k h_i + \sum_{j=1}^p \zeta_j^k \widehat{\mu}_j^k - \sum_{j=1}^p \widehat{\mu}_j^k g_j^- - \ell_k \sum_{j=1}^p \widehat{\mu}_j^k g_j^+ \\ &\leq 2\|\nabla f\| \|w^k\| + 2\|\nabla f\| \sum_{i=1}^m |\widehat{\lambda}_i^k h_i| - \sum_{j=1}^p \widehat{\mu}_j^k g_j^- + 2\|\nabla f\| \sum_{j=1}^p \widehat{\mu}_j^k g_j^+, \end{split}$$

where in the inequality we use $\zeta_j^k \widehat{\mu}_j^k \leq 0$ for all j. It follows from (10) that the right-hand side of the above inequality vanishes as k goes to infinity, and thus $(s^k, \ell_k, d^k) \to (0, 0, 0)$ by (11). This concludes the proof.

Theorem 3.1 indicates that to obtain CAKKT points, one can use the gradient projection $\operatorname{proj}_{\mathbb{L}^+_\Omega(x)}(0,0,-\nabla f(x))$ and try to find a mechanism that forces it to vanish across iterations. This is the essence of IR methods. Modern IR approaches that uses the traditional linearization (3), such as the one proposed in [16], employ projections that use the gradient of Lagrangian $\nabla_x \mathcal{L}(x,\lambda,\mu)$ instead of just $\nabla f(x)$, where the multipliers remain bounded during the execution of the algorithm. Since one can choose null multipliers as estimates, such strategies are more general. In order for our proposal to benefit from the use of the Lagrangian, we prove next that forcing the projection of $(0,0,-\nabla_x\mathcal{L})$ onto \mathbb{L}^+_Ω to vanish also guarantees CAKKT points. This extends the first part of Theorem 3.1.

Corollary 3.2. Let \bar{x} be a feasible point of (NLP). Then \bar{x} is a CAKKT point with associated sequence $\{x^k\}$ if, and only if, $\{x^k\} \to \bar{x}$ and there are bounded sequences $\{\bar{\lambda}^k\} \subset \mathbb{R}^m$, $\{\bar{\mu}^k\} \subset \mathbb{R}^p$ such that $\min\{-g(x^k), \bar{\mu}^k\} \to 0$ and

$$\operatorname{proj}_{\mathbb{L}_{\Omega}^{+}(x^{k})}(0,0,-\nabla_{x}\mathcal{L}(x^{k},\bar{\lambda}^{k},\bar{\mu}^{k})) \to 0.$$

Proof. If \bar{x} is a CAKKT point then, by Theorem 3.1, the projection of the Lagrangian vanishes taking $\bar{\lambda}^k = 0$ and $\bar{\mu}^k = 0$, $\forall k$, since $\nabla_x \mathcal{L}(x^k, 0, 0) = \nabla f(x^k)$.

Let us prove the converse. Following the proof of Theorem 3.1, but applied to the projection $\operatorname{proj}_{\mathbb{L}_{0}^{+}(x^{k})}(0,0,-\nabla_{x}\mathcal{L}(x^{k},\bar{\lambda}^{k},\bar{\mu}^{k}))$, we find sequences $\{\lambda^{k}\}\subset\mathbb{R}^{m}$, $\{\mu^{k}\}\subset\mathbb{R}^{p}_{+}$ such that

$$\nabla_x \mathcal{L}(x^k, \bar{\lambda}^k, \bar{\mu}^k) + \nabla h(x^k) \lambda^k + \nabla g(x^k) \mu^k \to 0, \quad \lambda_i^k h_i(x^k) \to 0, \forall i, \quad \mu_j^k g_j(x^k) \to 0, \forall j.$$
 (12)

We define $\widehat{\lambda}^k := \lambda^k + \overline{\lambda}^k$ and $\widehat{\mu}^k := \mu^k + \overline{\mu}^k$. From the feasibility of \overline{x} , the boundedness of $\{\overline{\lambda}^k\}$ and $\{\overline{\mu}^k\}$, and from $\overline{\mu}_j^k \to 0$ when $g_j(\overline{x}) < 0$, we have $\widehat{\lambda}_i^k h_i(x^k) \to 0$ for all i and $\widehat{\mu}_j^k g_j(x^k) \to 0$ for all j. Thus, the statement follows from (12).

Inspired in Corollary 3.2, we propose Algorithm 1 to solve (NLP). Note that its optimization step carries the projection of the Lagrangian onto $\mathbb{L}^+_{\Omega}(y^{k+1})$ and the multiplier estimates generated are bounded by Step 2.

In Section 5 we prove that the new IR method (Algorithm 1) converges to CAKKT points of (NLP) without exogenous assumptions. Actually, Algorithm 1 generates CAKKT sequences. We also mention that all necessary hypotheses for the new IR method are quite common in this context. To the best of our knowledge, our method is the first IR strategy that theoretically improves the convergence of the IR methodology, historically linked to the weaker AGP condition. Compared to previous IR methods, notably [16], this is done with a small modification in the optimization step.

We conclude this section by comparing the linearized sets (3) and (4). As we already mention, while $d \in \mathbb{L}_{\Omega}(x)$ implies $(0,0,d) \in \mathbb{L}_{\Omega}^+(x)$, the converse is not always true at infeasible points. Regarding Step 3 of Algorithm 1, this means that we have more freedom to generate non-null directions d than the classic IR methods, that use $\mathbb{L}_{\Omega}(x)$, and thus we have more chances of decreasing the objective function or the Lagrangian. Geometrically, the possibly larger set of directions defined by $\mathbb{L}_{\Omega}^+(x)$ coincides with that of $\mathbb{L}_{\Omega}(x)$ as $(s,\ell) \to 0$; see Figure 1. So, it is expected that IR methods with $\mathbb{L}_{\Omega}^+(x)$ have a more stringent stopping criteria than others, restricting the possible limit points of the method. This is illustrated by the convergence to CAKKT points in contrast with AGP of other IR strategies. Therefore, we believe that defining different linearizations is a fundamental step to obtain stronger IR methods.

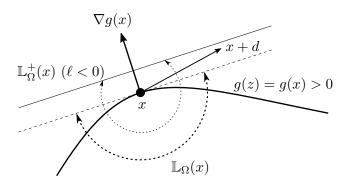


Figure 1: Comparison between $\mathbb{L}_{\Omega}(x)$ and $\mathbb{L}_{\Omega}^{+}(x)$. In the figure, $g(x) \leq 0$ is the unique constraint and x is infeasible (i.e., $g^{+}(x) = g(x) > 0$, $g^{-}(x) = 0$). The cone $\mathbb{L}_{\Omega}(x) = \{d \mid \nabla g(x)^{T} d \leq 0\}$ is a proper subset of $\{d \mid \ell g(x) + \nabla g(x)^{T} d \leq 0\}$ whenever $\ell < 0$, which is the "d-part" of $\mathbb{L}_{\Omega}^{+}(x)$. When $\ell \to 0^{-}$, these sets tend to be equal.

4 The new Inexact Restoration method

In this section, we describe our IR method. It is based on the general IR framework introduced in [20] but employs a novel optimization step. The convergence theory will be presented in the next section.

Consider the measure of infeasibility $\phi: \mathbb{R}^n \to \mathbb{R}_+$ defined as

$$\phi(x) = ||g^{+}(x)|| + ||h(x)||$$

and the merit function $\Phi: \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^p_+ \times [0,1] \to \mathbb{R}$ define as

$$\Phi(x,\lambda,\mu,\theta) := \theta \mathcal{L}(x,\lambda,\mu) + (1-\theta)\phi(x). \tag{13}$$

This merit function is the sharp Lagrangian used in previous IR methods, see [16, 29].

Our general IR framework is presented in Algorithm 1.

Now, some comments concerning Algorithm 1:

Remark 1. As in [16], we deal with the merit function (13), which is a convex combination of the Lagrangian and the infeasibility measure ϕ . The IR algorithm proposed in [16] does not have safeguards for the Lagrange multiplier estimates as in Step 2 of Algorithm 1, but they assume that the generated sequence of multipliers is bounded. Instead, we prefer to state our algorithm with safeguards, which is a natural way to bound multiplier estimates in practical implementations. In particular, such safeguards are used in popular implementations of AL methods [14]. In our framework, we do not explicitly use any update rules for the multipliers, we just require the general conditions in Step 2. In section 6 we discuss the update rule used in our implementation.

Remark 2. Originally, in [30], the IR method was proposed using the merit function $\theta f(x) + (1 - \theta)\phi(x)$ instead of (13). Therefore, no Lagrange multiplier is computed and the subproblem (14) is stated with $\nabla f(y^{k+1})$ instead of $\nabla_x \mathcal{L}(y^{k+1}, \lambda^{k+1}, \mu^{k+1})$. Of course, this case is recovered by taking $\lambda_{\min} = \lambda_{\max} = \mu_{\max} = 0$. The use of the Lagrangian was first proposed in [29]. Since then, the merit function (13) has been used in several works; see for example [16] and references therein. Also, it has been suggested that the use of the Lagrangian may reduce Maratos-like effects [17].

Remark 3. The last condition on μ_j^{k+1} in Step 2 meets the requirement in Corollary 3.2. This is natural, since Lagrange multipliers associated with inactive inequality constraints are zero. It is worth mentioning that in previous works on IR methods, such as [16, 29], only equality constraints are explicitly considered; inequalities are treated by inserting slack variables $z \geq 0$.

In Step 3 of Algorithm 1 we require that (s^k, ℓ_k, d^k) must be an ε_k -approximate solution of (14), in the following sense:

• it is almost feasible, i.e.,

$$|s_i^k h_i(y^{k+1}) + \nabla h_i(y^{k+1})^T d^k| \le \varepsilon_k, \quad \forall i,$$

$$g_j^-(y^{k+1}) + \ell_k g_j^+(y^{k+1}) + \nabla g_j(y^{k+1})^T d^k \le \varepsilon_k, \quad \forall j;$$
(17)

• it has an almost non-positive objective value, i.e.,

$$\nabla_x \mathcal{L}(y^{k+1}, \lambda^{k+1}, \mu^{k+1})^T d^k + \frac{1}{2} (d^k)^T H_k d^k + \frac{1}{2} (s^k)^T s^k + \frac{1}{2} \ell_k^2 \le \varepsilon_k.$$
 (18)

This is reasonable since (0,0,0) is feasible for (14), with null objective;

• it is almost optimal, i.e.,

$$\|\operatorname{proj}_{\mathbb{L}_{\mathcal{O}}^{+}(y^{k+1})}(-s^{k}, -\ell_{k}, -\nabla_{x}\mathcal{L}(y^{k+1}, \lambda^{k+1}, \mu^{k+1}) - H_{k}d^{k})\| \le \varepsilon_{k}.$$
 (19)

One interpretation of (19) is the following: at the solution $(\hat{s}, \hat{\ell}, \hat{d})$ of (14), the projection of $(-\hat{s}, -\hat{\ell}, -\nabla_x \mathcal{L}(y^{k+1}, \lambda^{k+1}, \mu^{k+1}) - H_k \hat{d})$ onto $\mathbb{L}^+_{\Omega}(y^{k+1})$ vanishes. So (19) says that (s^k, ℓ_k, d^k) is almost as good as $(\hat{s}, \hat{\ell}, \hat{d})$. Also, due to the convexity of the subproblem (as long as some positivity on H_k is imposed, see assumption H2 or H2' below), $(\hat{s}, \hat{\ell}, \hat{d})$ is a solution if and only if the projection vanishes.

Algorithm 1 General Inexact Restoration framework

Set $r \in [0,1)$, $\theta_{-1} \in (0,1)$, $-\infty < \lambda_{\min} \le \lambda_{\max} < \infty$, $\mu_{\max} \ge 0$, $\sigma_{\min} > 0$ and $\gamma_c \in (0,1/2]$. Choose a non-negative summable sequence $\{\varepsilon_k\}$, and $x^0 \in \mathbb{R}^n$, $\lambda^0 \in [\lambda_{\min}, \lambda_{\max}]^m$ and $\mu^0 \in [0, \mu_{\max}]^p$. Set $k \leftarrow 0$.

- 1. Restoration step. If $\phi(x^k) = 0$, define $y^{k+1} := x^k$. If $\phi(x^k) > 0$, compute, if possible, a point y^{k+1} such that $\phi(y^{k+1}) \le r\phi(x^k)$. Otherwise, abort the execution declaring failure.
- 2. Estimation of Lagrange multipliers. Compute $\lambda^{k+1} \in [\lambda_{\min}, \lambda_{\max}]^m$ and $\mu^{k+1} \in [0, \mu_{\max}]^p$. The sequence $\{\mu^{k+1}\}$ must be chosen so that $\min\{-g(y^{k+1}), \mu^{k+1}\} \to 0$.
- 3. Optimization step. Compute a symmetric matrix H_k and find an ε_k -approximate solution (defined below, see (17)–(19)) (s^k, ℓ_k, d^k) of the quadratic problem

$$\min \nabla_{x} \mathcal{L}(y^{k+1}, \lambda^{k+1}, \mu^{k+1})^{T} d + \frac{1}{2} d^{T} H_{k} d + \frac{1}{2} s^{T} s + \frac{1}{2} \ell^{2}$$

$$s_{i} h_{i}(y^{k+1}) + \nabla h_{i}(y^{k+1})^{T} d = 0, \quad \forall i,$$
s.t.
$$g_{j}^{-}(y^{k+1}) + \ell g_{j}^{+}(y^{k+1}) + \nabla g_{j}(y^{k+1})^{T} d \leq 0, \quad \forall j.$$

$$(14)$$

4. Penalty parameter computation. If

$$\Phi(y^{k+1}, \lambda^{k+1}, \mu^{k+1}, \theta_{k-1}) \le \Phi(x^k, \lambda^k, \mu^k, \theta_{k-1}) - \frac{1-r}{2} (\phi(x^k) - \phi(y^{k+1})),$$

set $\theta_k := \theta_{k-1}$. Otherwise, compute

$$\theta_k := \frac{(1+r)(\phi(x^k) - \phi(y^{k+1}))}{2(\mathcal{L}(y^{k+1}, \lambda^{k+1}, \mu^{k+1}) - \mathcal{L}(x^k, \lambda^k, \mu^k) + \phi(x^k) - \phi(y^{k+1}))}.$$

5. Globalization. If $d^k = 0$, set $t_k := 1$ and go to Step 6. Otherwise, compute by some backtracking strategy $t_k \in [0,1]$ such that

$$\mathcal{L}(y^{k+1} + t_k d^k, \lambda^{k+1}, \mu^{k+1})$$

$$\leq \mathcal{L}(y^{k+1}, \lambda^{k+1}, \mu^{k+1}) - \frac{\sigma_{\min}}{4} t_k \|d^k\|^2 - \gamma_c t_k \|(s^k, \ell_k)\|^2 + t_k \varepsilon_k$$
(15)

and

$$\Phi(y^{k+1} + t_k d^k, \lambda^{k+1}, \mu^{k+1}, \theta_k)
\leq \Phi(x^k, \lambda^k, \mu^k, \theta_k) - \frac{1-r}{4} (\phi(x^k) - \phi(y^{k+1})) + 2t_k \|\mathbf{1}\| \varepsilon_k. \tag{16}$$

Moreover, the backtracking procedure should be such that either $t_k = 1$ or there exists $\hat{t}_k \in [t_k, at_k]$ (a > 1) such that (15) or (16) fail.

6. Set $x^{k+1} := y^{k+1} + t_k d^k$. Update $k \leftarrow k+1$ and go to Step 1.

Of course, if we solve (14) "exactly" then all three conditions above hold with $\varepsilon_k = 0$. Thus, we could establish our theory just for this case, as done in previous works, e.g. [16]. However, we consider $\varepsilon_k \geq 0$ because (i) this opens up the possibility of using iterative methods to solve

(14), and (ii) the additional work needed to cover this case is negligible. It is worth mentioning that a similar IR algorithm, allowing inexact resolution of the optimization step problem, was considered in [18].

4.1 Well-definiteness of the method

The well-definiteness of Algorithm 1 is established under the following assumptions:

- **H1** ∇f , ∇h and ∇g are Lipschitz with constants L_F , L_H and L_G , respectively.
- **H2** In Step 3 of Algorithm 1, the matrices H_k are chosen such that their eigenvalues are uniformly in $[\sigma_{\min}, \sigma_{\max}]$, where $\sigma_{\max} \geq \sigma_{\min}$.

Hypothesis H2 ensures that (14) is strictly convex, having a unique and isolated minimizer. There are practical ways to achieve H2, for example, computing the Cholesky factorization of $H_k = \nabla_{xx}^2 \mathcal{L}(y^{k+1}, \lambda^{k+1}, \mu^{k+1}) + \sigma I$ by successively increasing $\sigma \geq 0$. On the other hand, (14) is a strictly convex problem when the Hessian of its objective function,

$$\hat{H}_k = \left[\begin{array}{cc} H_k & 0 \\ 0 & I_{m+1} \end{array} \right],$$

satisfies $\omega^T \hat{H}_k \omega \ge \sigma_{\min} \|\omega\|^2$ for all feasible $\omega = (d, s, \ell)$. In this sense, we enunciate the following relaxed version of H2:

H2' For $k \in \mathbb{N}$, the matrix H_k from Step 3 of Algorithm 1 is chosen such that $(d^k)^T H_k d^k \ge \sigma_{\min} \|d^k\|^2$, where d^k is the d-part of the computed solution, and its largest eigenvalue does not exceed $\sigma_{\max} \ge \sigma_{\min}$.

Our theory is valid with the above condition in place of H2 (especially Theorem 4.3 and Lemmas 4.2, 5.3 and 5.4), but it is not always implementable since d^k depends on H_k . For this reason, we prefer to state our results using H2. However, it is possible to achieve H2' in certain situations, such as when solving (14) "exactly" using an active-set method, whose equality-constrained quadratic problems are solved through its KKT optimality system after correcting inertia; see [33]. This approach is adopted in our implementation. It is worth mentioning that the convexity of the subproblem (14), guaranteed by assumption H2 or H2', typically requires a line search strategy like that in step 5 of Algorithm 1. However, this could be avoided in globalization techniques using trust region strategies as done in [9, 29, 30], which could favor fast local convergence.

The next theorem establishes the well-definiteness of Algorithm 1. Its proof is analogous to that of [16, Theorem 2.1] taking into account Lemmas 4.1 and 4.2. We remark that in [16], the conclusions of these lemmas are embedded in the method, and the authors use a quadratic problem related to (14) to show how such conditions can be achieved. Instead, we believe that putting (14) directly in the algorithm turns it clearer. We start with Lemma 4.1, which bounds the growth of infeasibility ϕ along the direction d^k computed in the optimization phase.

Lemma 4.1. If H1 holds then there exists $\gamma_{\phi} > 0$, which is independent of k, such that, for every $t \in [0,1]$,

$$\phi(y^{k+1} + td^k) \le \phi(y^{k+1}) + t(|\ell_k| + ||s^k||)\phi(y^{k+1}) + \gamma_\phi t^2 ||d^k||^2 + 2t||\mathbf{1}||\varepsilon_k.$$

Proof. First let us prove that

$$||g^{+}(y^{k+1} + td^{k})|| \le ||g^{+}(y^{k+1})|| + t|\ell_{k}|||g^{+}(y^{k+1})|| + \gamma_{q}t^{2}||d^{k}||^{2} + t||\mathbf{1}||\varepsilon_{k}$$
(20)

for every $t \in [0,1]$ and some $\gamma_g > 0$. For simplicity, we will omit the indices k and k+1. By Taylor's formula and Lipschitz continuity of ∇g , for all $t \geq 0$, we have

$$||g(y+td) - g(y) - t\nabla g(y)^T d|| = t \left\| \int_0^1 (\nabla g(y+tsd)^T d - \nabla g(y)^T d) ds \right\|$$

$$\leq t \int_0^1 L_G ||y+tsd-y|| ||d|| ds \leq \gamma_g t^2 ||d||^2,$$

where $\gamma_q := (1/2)L_G$. From (17), we have, for each $t \in [0,1]$,

$$g^{-}(y) + t\ell g^{+}(y) + t\nabla g(y)^{T}d - t\mathbf{1}\varepsilon = (1-t)g^{-}(y) + t[g^{-}(y) + \ell g^{+}(y) + \nabla g(y)^{T}d - \mathbf{1}\varepsilon] \le 0.$$

If $g_j(y+td) \leq 0$ then $g_j^+(y+td) = 0 \leq |g_j(y+td) - z_j|$ for all $z_j \in \mathbb{R}$. If $g_j(y+td) > 0$ then $g_j^+(y+td) = g_j(y+td) \leq g_j(y+td) - z_j$ for all $z_j \leq 0$. Thus

$$||g^{+}(y+td)|| \le ||g(y+td) - z||, \quad \forall z \in \mathbb{R}_{-}^{p}.$$

If we choose $z = g^-(y) + t\ell g^+(y) + t\nabla g(y)^T d - t\mathbf{1}\varepsilon$, we arrive at

$$||g^{+}(y+td)|| \leq ||g(y+td) - (g^{-}(y) + \ell t g^{+}(y) + t \nabla g(y)^{T} d - t \mathbf{1}\varepsilon)||$$

$$\leq ||g(y+td) - g(y) - t \nabla g(y)^{T} d|| + ||g(y) - g^{-}(y) - t \ell g^{+}(y)|| + ||t \mathbf{1}\varepsilon||$$

$$\leq \gamma_{g} t^{2} ||d||^{2} + |1 - t\ell|||g^{+}(y)|| + t||\mathbf{1}||\varepsilon$$

$$\leq \gamma_{g} t^{2} ||d||^{2} + ||g^{+}(y)|| + t||\ell|||g^{+}(y)|| + t||\mathbf{1}||\varepsilon,$$

where in the third inequality we use the equality $g(y) - g^{-}(y) = g^{+}(y)$.

Analogously, we can prove that

$$||h(y^{k+1} + td^k)|| \le ||h(y^{k+1})|| + t||s^k|| ||h(y^{k+1})|| + \gamma_h t^2 ||d^k||^2 + t||\mathbf{1}||\varepsilon_k$$
 (21)

for every $t \in [0,1]$ and some $\gamma_h > 0$. The statement follows from (20)–(21) with $\gamma_{\phi} = \gamma_g + \gamma_h$. Note that γ_{ϕ} does not depend on k.

The next lemma says that d^k is a direction in which $\mathcal{L}(\cdot, \lambda^{k+1}, \mu^{k+1})$ decreases locally from y^{k+1} or, at least, does not grow too much. Note that H1 implies the Lipschitz continuity of the gradient of $\mathcal{L}(\cdot, \lambda^{k+1}, \mu^{k+1})$ with a constant $L_{\mathcal{L}}$ that is independent of k (here, we use the fact that $\{\lambda^k\}$ are $\{\mu^k\}$ are bounded sequences).

Lemma 4.2. Assume that H1 and H2 hold. Then there exists $\bar{t} \in (0,1]$, which is independent of k, such that, for every $t \in [0,\bar{t}]$ and $\gamma_c \in (0,1/2]$,

$$\mathcal{L}(y^{k+1} + td^k, \lambda^{k+1}, \mu^{k+1}) \leq \mathcal{L}(y^{k+1}, \lambda^{k+1}, \mu^{k+1}) - \frac{\sigma_{\min}}{4} t \|d^k\|^2 - \gamma_c t \|(s^k, \ell_k)\|^2 + t\varepsilon_k.$$

Proof. For simplicity, we will omit the indices k and k+1. From (18) and H2,

$$2\nabla_x \mathcal{L}(y,\lambda,\mu)^T d \le -d^T H d - s^T s - \ell^2 + \varepsilon \le -\sigma_{\min} \|d\|^2 - \|s\|^2 - \ell^2 + 2\varepsilon.$$

So, by Taylor's formula of $x \to \mathcal{L}(x, \lambda^k, \mu^k)$ at y and the Lipschitz continuity of its gradient, we have

$$\mathcal{L}(y + td, \lambda, \mu) \leq \mathcal{L}(y, \lambda, \mu) + t\nabla_{x}\mathcal{L}(y, \lambda, \mu)^{T}d + (1/2)L_{\mathcal{L}}t^{2}\|d\|^{2}$$

$$\leq \mathcal{L}(y, \lambda, \mu) + (1/2)t\left[-\sigma_{\min}\|d\|^{2} - \|s\|^{2} - \ell^{2} + 2\varepsilon\right] + (1/2)L_{\mathcal{L}}t^{2}\|d\|^{2}$$

$$= \mathcal{L}(y, \lambda, \mu) - (1/2)t(\sigma_{\min} - L_{\mathcal{L}}t)\|d\|^{2} - (1/2)t\|(s, \ell)\|^{2} + t\varepsilon$$

$$\leq \mathcal{L}(y, \lambda, \mu) - (\sigma_{\min}/4)t\|d\|^{2} - \gamma_{c}t\|(s, \ell)\|^{2} + t\varepsilon$$

for every $t \in [0, \bar{t}]$ and $\gamma_c \in (0, 1/2]$, where $\bar{t} = \min\{1, \sigma_{\min}/(2L_{\mathcal{L}})\}$.

Now, we present the main result of this section.

Theorem 4.3. Assume valid H1 and H2. If y^{k+1} is successfully computed in the restoration step of Algorithm 1, then x^{k+1} is well-defined.

Proof. Clearly, Step 2 is accomplished by taking zero multipliers, although there are other possibilities (see section 6.1). By H2, (14) always has a solution, so Step 3 is well-defined. If the inequality in Step 4 does not hold, then after straightforward calculations we arrive at

$$\theta_{k-1} \Big[\mathcal{L}(y^{k+1}, \lambda^{k+1}, \mu^{k+1}) - \mathcal{L}(x^k, \lambda^k, \mu^k) + \phi(x^k) - \phi(y^{k+1}) \Big]$$

$$> \frac{1+r}{2} \Big(\phi(x^k) - \phi(y^{k+1}) \Big) \ge 0.$$

Therefore, Step 4 is well-defined. Lemma 4.2 guarantees that condition (15) holds for all t > 0 small enough. It remains to prove that the same occurs with (16). In this case, the requirement over the accepted step-size t_k ($t_k = 1$ or $\hat{t}_k \in [t_k, at_k]$) can be achieved by a simple backtracking procedure, see section 6.1.4.

Let y^{k+1} be a point computed in Step 1. By the way that θ_k is chosen in Step 4, we have

$$\Phi(y^{k+1}, \lambda^{k+1}, \mu^{k+1}, \theta_k) \le \Phi(x^k, \lambda^k, \mu^k, \theta_k) - \frac{1-r}{2} (\phi(x^k) - \phi(y^{k+1})). \tag{22}$$

We split the proof into two cases: first, suppose that $\phi(x^k) > \phi(y^{k+1})$. By (22) and the continuity of Φ , we have $\Phi(y^{k+1} + td^k, \lambda^{k+1}, \mu^{k+1}, \theta_k) \leq \Phi(x^k, \lambda^k, \mu^k, \theta_k) - \frac{1-r}{4}(\phi(x^k) - \phi(y^{k+1})) + 2t\|\mathbf{1}\|\varepsilon_k$ for all t > 0 small enough. This implies the well definiteness of the backtracking procedure in Step 5, and the statement follows.

Now, suppose that $\phi(x^k) = \phi(y^{k+1})$. From Step 1, $\phi(x^k) = \phi(y^{k+1}) = 0$ and $y^{k+1} = x^k$, and thus (22) implies $\Phi(y^{k+1}, \lambda^{k+1}, \mu^{k+1}, \theta_k) \leq \Phi(x^k, \lambda^k, \mu^k, \theta_k)$. Furthermore, from Lemma 4.1 we have $\phi(y^{k+1} + td^k) \leq \gamma_{\phi}t^2 ||d^k||^2 + 2t||\mathbf{1}||\varepsilon_k$ for every $t \leq 1$. Thus, from the previous inequalities, the definition of Φ and Lemma 4.2 we obtain, after straightforward calculations,

$$\begin{split} \Phi(y^{k+1} + td^k, \lambda^{k+1}, \mu^{k+1}, \theta_k) &\leq \Phi(x^k, \lambda^k, \mu^k, \theta_k) \\ -t\|d^k\|^2 [\theta_k(\sigma_{\min}/4) - (1 - \theta_k)\gamma_\phi t] - \theta_k \gamma_c t\|(s^k, \ell_k)\|^2 + t\varepsilon_k [\theta_k + 2(1 - \theta_k)\|\mathbf{1}\|] \end{split}$$

for all t > 0 sufficiently small. The statement follows from the above inequality, noting that $t\varepsilon_k[\theta_k + 2(1-\theta_k)\|\mathbf{1}\|] \le 2t\|\mathbf{1}\|\varepsilon_k$. This concludes the proof.

5 Convergence of the proposed IR method

Now, we analyze the convergence of Algorithm 1. Throughout this section, we assume that the IR method does not stop after a finite number of iterations. We start with the following hypotheses:

H3 All the iterates y^k stay in a bounded set.

H4 There is $\beta > 0$ such that $\mathcal{L}(y^{k+1}, \lambda^{k+1}, \mu^{k+1}) - \mathcal{L}(x^k, \lambda^k, \mu^k) \leq \beta \phi(x^k)$ for all k.

Assumption H3 is satisfied, for example, if $\{x^k\}$ remains bounded and if we compute y^{k+1} within a neighborhood of x^k in Step 1 (note that, due to the continuity of g, such a valid y^{k+1} always exists). This can be done by applying an iterative method (or only some steps of that) for minimizing $\phi(x)$ subject to such a neighborhood. On the other hand, $\{x^k\}$ remains bounded if the original problem (NLP) has box-constraints, as example. Also, it can be expected that, maintaining y^{k+1} near x^k , assumption H4 holds whenever x^k is infeasible, at least when the multipliers do not change too much. A similar assumption was used in [16]. However, to simplify the exposition, we have decided not to impose any additional explicit algorithmic condition on the point y^{k+1} at Step 1 of Algorithm 1.

Assumption H4 is standard in the literature of IR methods, see for instance [16, 20]. In some sense, this condition limits the increase of the objective function of (14) that might be incurred by moving towards feasibility. It is satisfied under some CQ, see [20, Remark 3]. Recently, a procedure to obtain H4 through a numerical method in the restoration phase was provided [18]. To simplify our presentation, we do not follow this approach; instead, we explicitly assume the validity of H4.

5.1 Technical results

In the next lemmas, we exhibit some useful properties mainly derived from solving the problem in Step 3 of Algorithm 1. Sequences $\{x^k\}$, $\{y^{k+1}\}$, $\{\lambda^{k+1}\}$, $\{\mu^{k+1}\}$, $\{s^k\}$, $\{\ell_k\}$ and $\{d^k\}$ are those generated by Algorithm 1.

Lemma 5.1. Under H4, the sequence $\{\theta_k\}$ is non-increasing and bounded away from zero, that is, there exists $\bar{\theta} \in (0,1)$ such that $\bar{\theta} \leq \theta_{k+1} \leq \theta_k$ for all k.

Proof. The proof is analogous to that of [16, Lemma 3.1].

Next, we prove that the line-search procedure in Step 5 of Algorithm 1 does not produce excessively short steps.

Lemma 5.2. Under H1–H4 we have $t_k \ge \underline{t}_k$ for all k, where

$$\underline{t}_k := \frac{1}{a} \min \left\{ \overline{t}, \frac{\overline{\theta} \gamma_d}{(1 - \overline{\theta}) \gamma_\phi}, \frac{(1 - r)^2}{4r(1 - \overline{\theta})(|\ell_k| + ||s^k||)} \right\},$$

 $\gamma_d = \sigma_{\min}/4$ and \bar{t} , γ_{ϕ} , $\bar{\theta}$ are given by Lemmas 4.1, 4.2 and 5.1.

Proof. Let us fix the iteration index k. As $\overline{t} \leq 1$ and a > 1, we have $\underline{t}_k \leq a\underline{t}_k \leq 1$. Let d^k be the solution of (14). If $d^k = 0$, we have $t_k = 1 \geq \underline{t}_k$ by Step 5 of Algorithm 1, and then there is nothing left to prove.

Now, consider the case $d^k \neq 0$. For simplicity, denote $\nu^j := (\lambda^j, \mu^j, \theta_k)$ for j = k, k + 1. Let us prove that (16) holds for every $0 < t \le a\underline{t}_k$. From Step 4, θ_k is defined so that

$$\Phi(y^{k+1}, \nu^{k+1}) \le \Phi(x^k, \nu^k) - \frac{1-r}{2} \Big(\phi(x^k) - \phi(y^{k+1}) \Big).$$

Adding and subtracting $\Phi(y^{k+1} + td^k, \nu^{k+1})$ to the above inequality, we arrive at

$$\Phi(y^{k+1} + td^k, \nu^{k+1}) - \Phi(x^k, \nu^k)
\leq \Phi(y^{k+1} + td^k, \nu^{k+1}) - \Phi(y^{k+1}, \nu^{k+1}) - \frac{1-r}{2} \Big(\phi(x^k) - \phi(y^{k+1}) \Big).$$
(23)

On the other hand, we have

$$\Phi(y^{k+1} + td^{k}, \nu^{k+1}) - \Phi(y^{k+1}, \nu^{k+1})
= \theta \left[\mathcal{L}(y + td, \lambda, \mu) - \mathcal{L}(y, \lambda, \mu) \right] + (1 - \theta) \left[\phi(y + td) - \phi(y) \right]
\leq \theta \left[t\varepsilon - \gamma_{d}t \|d\|^{2} \right] + (1 - \theta) \left[t(|\ell| + \|s\|)\phi(y) + \gamma_{\phi}t^{2} \|d\|^{2} + 2t \|\mathbf{1}\|\varepsilon \right]
= 2t \|\mathbf{1}\|\varepsilon - t\theta\varepsilon(2\|\mathbf{1}\| - 1) + t(1 - \theta)(|\ell| + \|s\|)\phi(y) - t(\theta\gamma_{d} - (1 - \theta)\gamma_{\phi}t) \|d\|^{2}
\leq 2t \|\mathbf{1}\|\varepsilon + t(1 - \theta)(|\ell| + \|s\|)\phi(y) - t(\theta\gamma_{d} - (1 - \theta)\gamma_{\phi}t) \|d\|^{2},$$
(24)

for any $t \in [0, \bar{t}]$, where the first inequality follows from Lemmas 4.1 and 4.2 (we omit indices k and k+1 for simplicity). The last term on the right-hand side of (24) is non-negative if $\theta_k \gamma_d - (1 - \theta_k) \gamma_\phi t \leq 0$, in particular, as $\bar{\theta} \leq \theta_k$ by Lemma 5.1, for all

$$t \le \tilde{t}_1 := \min \left\{ \bar{t} , \frac{\bar{\theta} \gamma_d}{(1 - \bar{\theta}) \gamma_\phi} \right\} \le \frac{\theta_k \gamma_d}{(1 - \theta_k) \gamma_\phi}.$$

Thus, (24) and $\phi(y^{k+1}) \leq r\phi(x^k)$ (see Step 1 of Algorithm 1) imply

$$\Phi(y^{k+1} + td^k, \nu^{k+1}) - \Phi(y^{k+1}, \nu^{k+1}) - 2t \|\mathbf{1}\| \varepsilon_k - \frac{1-r}{4} (\phi(x^k) - \phi(y^{k+1}))
\leq t(1-\theta_k)(|\ell_k| + \|s^k\|)\phi(y^{k+1}) - \frac{1-r}{4} (\phi(x^k) - \phi(y^{k+1}))
\leq t(1-\theta_k)(|\ell_k| + \|s^k\|)r\phi(x^k) - \frac{1-r}{4} (\phi(x^k) - r\phi(x^k))
= \left(t(1-\theta_k)(|\ell_k| + \|s^k\|)r - \frac{(1-r)^2}{4}\right)\phi(x^k)$$
(25)

for all $t \leq \tilde{t}_1$. If $|\ell_k| + ||s^k|| = 0$, the right side of (25) is non-positive for all t; otherwise, i.e. $|\ell_k| + ||s^k|| > 0$, it is non-positive for all

$$t \le \tilde{t}_2 := \frac{(1-r)^2}{4r(1-\bar{\theta})(|\ell_k| + ||s^k||)} \le \frac{(1-r)^2}{4r(1-\theta_k)(|\ell_k| + ||s^k||)}.$$

For convenience, define $\tilde{t}_2 := \infty$ if $|\ell_k| + ||s^k|| = 0$. In any case, (23) and (25) give

$$\Phi(y^{k+1} + td^k, \nu^{k+1}) - \Phi(x^k, \nu^k) \le -\frac{1-r}{4} \left(\phi(x^k) - \phi(y^{k+1}) \right) + 2t \|\mathbf{1}\| \varepsilon_k$$

for all $t \leq \min\{\tilde{t}_1, \tilde{t}_2\}$, and hence (16) holds for all $t \leq \min\{\tilde{t}_1, \tilde{t}_2\}$.

Finally, from the backtracking procedure we have

$$t_k \ge \frac{1}{a}\min\{\tilde{t}_1, \tilde{t}_2\} = \underline{t}_k.$$

This concludes the proof.

The next result shows that under H2 and H3, the sequence of ε_k -approximate solutions of (14) is bounded.

Lemma 5.3. Suppose valid H2 and H3. Then $\{(s^k, \ell_k, d^k)\}$ is a bounded sequence.

Proof. First notice that $\{\nabla_x \mathcal{L}(y^{k+1}, \lambda^{k+1}, \mu^{k+1})\}$ is bounded due to H3, the boundedness of $\{\lambda^{k+1}\}$ and $\{\mu^{k+1}\}$, and the continuity of all gradients. We split the proof into two cases, depending on whether $\{d^k\}$ is bounded or not. If it is bounded, (18) implies $||s^k||^2 + \ell_k^2 \leq 2||\nabla_x \mathcal{L}(y^{k+1}, \lambda^{k+1}, \mu^{k+1})|| ||d^k|| - (d^k)^T H_k d_k + 2\varepsilon_k$, and thus the sequence $\{(s^k, \ell_k)\}$ is bounded. Now, suppose that $\{d^k\}$ is unbounded. Without loss of generality, we assume that $||d^k|| \to \infty$ after taking a suitable subsequence. From (18) and H2, we get $2\nabla_x \mathcal{L}(y^{k+1}, \lambda^{k+1}, \mu^{k+1})^T d^k + \sigma_{\min} ||d^k||^2 \leq 2\varepsilon_k$. Dividing this expression by $||d^k||^2$ and passing to the limit we conclude that $\sigma_{\min} \leq 0$, contradicting H2. Thus, $\{d^k\}$ is bounded and so is $\{(s^k, \ell_k)\}$ as well.

The next result complements Lemma 5.3. It reveals how $\|(s^k, \ell_k, d^k)\|$ controls the gradient projection associated with CAKKT, which is a crucial fact for the main convergence theorem of the next section.

Lemma 5.4. Suppose valid H2. For every $\sigma \ge \max\{1, \sigma_{\max}\}$ and every k, we have $\|\operatorname{proj}_{\mathbb{L}^+_{\Omega}(y^{k+1})}(0, 0, -\nabla_x \mathcal{L}(y^{k+1}, \lambda^{k+1}, \mu^{k+1})\| \le \sigma \|(s^k, \ell_k, d^k)\| + \varepsilon_k$.

Proof. To simplify the presentation, we omit indices k and k+1. Let us define $(\widehat{s}, \widehat{\ell}, \widehat{d}) := \operatorname{proj}_{\mathbb{L}^+_{\Omega}(y)}(-s, -\ell, -\nabla_x \mathcal{L}(y, \lambda, \mu) - Hd)$. From (19) we have $\|(\widehat{s}, \widehat{\ell}, \widehat{d})\| \leq \varepsilon$. Thus, from the non-expansiveness of the projection,

$$\begin{aligned} &\|\operatorname{proj}_{\mathbb{L}_{\Omega}^{+}(y)}(0,0,-\nabla_{x}\mathcal{L}(y,\lambda,\mu))\| \\ &= \|\operatorname{proj}_{\mathbb{L}_{\Omega}^{+}(y)}(0,0,-\nabla_{x}\mathcal{L}(y,\lambda,\mu)) - \operatorname{proj}_{\mathbb{L}_{\Omega}^{+}(y)}(-s,-\ell,-\nabla_{x}\mathcal{L}(y,\lambda,\mu) - Hd) + (\widehat{s},\widehat{\ell},\widehat{d})\| \\ &\leq \|(s,\ell,-\nabla_{x}\mathcal{L}(y,\lambda,\mu)) + \nabla_{x}\mathcal{L}(y,\lambda,\mu)) + Hd)\| + \|(\widehat{s},\widehat{\ell},\widehat{d})\| \\ &= \|(s,\ell,Hd)\| + \|(\widehat{s},\widehat{\ell},\widehat{d})\| \leq \widehat{\sigma}\|(s,\ell,d)\| + \varepsilon, \end{aligned}$$

where $\hat{\sigma}$ is the norm of the linear operator $(s, \ell, d) \to (s, \ell, Hd)$. Finally, in view of H2 we have $\hat{\sigma} \leq \max\{1, \sigma_{\max}\}$, from which follows the statement.

Observe that Lemma 5.3 ensures the boundedness of the sequence $\{x^k\}$. In fact, since $x^{k+1} := y^{k+1} + t_k d^k$ (Step 6 of Algorithm 1), we have $\|x^{k+1}\| \le \|y^{k+1}\| + t_k \|d^k\|$ and the statement follows from Lemma 5.3, H3 and the fact that $t_k \le 1$. In particular, the sequences $\{f(x^k)\}$, $\{\phi(x^k)\}$ and $\{\mathcal{L}(x^k, \lambda^k, \mu^k)\}$ are bounded. Furthermore, every limit point of the IR method is feasible as the next lemma shows.

Lemma 5.5. Under H2, H3 and H4, the sequence $\{\phi(x^k)\}$ is summable.

Proof. We follow arguments similar to the proof of [16, Lemma 3.2]. From Steps 1 and 5 of Algorithm 1 and the fact that $t_k \leq 1$, we have, for all k,

$$\Phi(z^{k+1}, \theta_k) \le \Phi(z^k, \theta_k) - \frac{(1-r)^2}{4}\phi(x^k) + 2\|\mathbf{1}\|\varepsilon_k,$$

where $z^k := (x^k, \lambda^k, \mu^k)$. Using the definition of Φ , and after some straightforward computations, the above inequality implies

$$\frac{(1-r)^2}{4}\phi(x^k) \le \theta_k(\mathcal{L}(z^k) - \mathcal{L}(z^{k+1})) + (1-\theta_k)(\phi(x^k) - \phi(x^{k+1})) + 2\|\mathbf{1}\|\varepsilon_k.$$

By Lemma 5.1, $0 < \bar{\theta} \le \theta_k \le \theta_0$ for all k. So, dividing the above inequality by θ_k ,

$$\frac{(1-r)^2}{4\theta_0}\phi(x^k) \le \mathcal{L}(z^k) - \mathcal{L}(z^{k+1}) + \frac{1-\theta_k}{\theta_k}(\phi(x^k) - \phi(x^{k+1})) + 2\bar{\theta}^{-1} ||\mathbf{1}|| \varepsilon_k.$$
 (26)

The statement follows from (26) if the sequence $\{a_k(\phi(x^k) - \phi(x^{k+1}))\}$ is summable, where $a_k := (1 - \theta_k)/\theta_k$. From Lemma 5.1, it is straightforward to verify that $0 < (1 - \theta_0)/\theta_0 \le a_{k-1} \le a_k \le (1 - \bar{\theta})/\bar{\theta}$ for all k. Also,

$$a_k(\phi(x^k) - \phi(x^{k+1})) = (a_k - a_{k-1})\phi(x^k) + a_{k-1}\phi(x^k) - a_k\phi(x^{k+1}).$$
(27)

If $M := \sup_{k} \phi(x^k)$ then $\sum_{k=1}^{K} (a_k - a_{k-1})\phi(x^k) \leq \sum_{k=1}^{K} (a_k - a_{k-1})M = (a_K - a_0)M$. Thus, by summing (27) over k = 1, ..., K and noting that $a_K \phi(x^{K+1}) \geq 0$, we obtain

$$\sum_{k=1}^{K} a_k(\phi(x^k) - \phi(x^{k+1})) \le (a_K - a_0)M + a_0\phi(x^1) - a_K\phi(x^{K+1}) \le \frac{(1 - \bar{\theta})}{\bar{\theta}}M.$$

By summing (26) over k = 1, ..., K and using the last inequality, we arrive at

$$\frac{(1-r)^2}{4\theta_1} \sum_{k=1}^K \phi(x^k) \le \mathcal{L}(z^1) - \mathcal{L}(z^{K+1}) + \frac{(1-\bar{\theta})}{\bar{\theta}} M + 2\bar{\theta}^{-1} ||\mathbf{1}|| \sum_{k=1}^K \varepsilon_k.$$

Thus, $\{\phi(x^k)\}$ is summable because $\{\varepsilon_k\}$ is summable and $\{\mathcal{L}(z^k)\}$ is bounded.

Lemma 5.6. Under H1-H4, $\{t_k\}$ is bounded below by some $\underline{t} > 0$.

Proof. By Lemma 5.3, the sequence $\{(s^k, \ell_k, d^k)\}$ is bounded. Take M > 0 such that $||s^k|| + \ell_k \le M$ for all k. By Lemma 5.2, we have

$$t_k \ge \underline{t}_k \ge \underline{t} := \frac{1}{a} \min \left\{ \overline{t}, \frac{\overline{\theta} \gamma_d}{(1 - \overline{\theta}) \gamma_\phi}, \frac{(1 - r)^2}{4r(1 - \overline{\theta})M} \right\} > 0$$

for all k, concluding the proof.

Lemma 5.7. Suppose valid H1–H4, and let $\{x^k\}$ be an infinite sequence generated by Algorithm 1. Then $(s^k, \ell_k, d^k) \to 0$.

Proof. Let us denote $\nu^k := (\lambda^k, \mu^k)$. From H4 and Lemmas 4.2 and 5.6, we have

$$\begin{split} \gamma_{d}\underline{t}\|d^{k}\|^{2} + \gamma_{c}\underline{t}\|(s^{k}, \ell_{k})\|^{2} &\leq \mathcal{L}(y^{k+1}, \nu^{k+1}) - \mathcal{L}(y^{k+1} + t_{k}d^{k}, \nu^{k+1}) + t_{k}\varepsilon_{k} \\ &= \left[\mathcal{L}(y^{k+1}, \nu^{k+1}) - \mathcal{L}(x^{k}, \nu^{k})\right] + \mathcal{L}(x^{k}, \nu^{k}) - \mathcal{L}(x^{k+1}, \nu^{k+1}) + t_{k}\varepsilon_{k} \\ &\leq \beta\phi(x^{k}) + \mathcal{L}(x^{k}, \nu^{k}) - \mathcal{L}(x^{k+1}, \nu^{k+1}) + t_{k}\varepsilon_{k} \end{split}$$

for all k, where $\gamma_d = \sigma_{\min}/4$. By summing over k = 1, ..., K and using the fact that $t_k \leq 1$, we obtain

$$\sum_{k=1}^{K} (\gamma_{d}\underline{t} \|d^{k}\|^{2} + \gamma_{c}\underline{t} \|(s^{k}, \ell_{k}^{2})\|^{2}) \leq \mathcal{L}(x^{1}, \nu^{1}) - \mathcal{L}(x^{K+1}, \nu^{K+1}) + \sum_{k=1}^{K} (\beta \phi(x^{k}) + \varepsilon_{k}).$$

Taking $K \to \infty$, we conclude that $\{\|d^k\|^2\}$, $\{\|s^k\|^2\}$ and $\{|\ell_k|^2\}$ are summable using the definition of $\{\varepsilon_k\}$ in Algorithm 1, the boundedness of $\{(x^k, \lambda^k, \mu^k)\}$ and Lemma 5.5.

5.2 The main convergence result

Next, we state the main convergence result for our IR framework, which is supported by the previous lemmas.

Theorem 5.8. Suppose that H1–H4 hold and consider the infinite sequences $\{x^k\}$ and $\{y^k\}$ generated by Algorithm 1. Then every limit point of $\{x^k\}$ (or $\{y^k\}$) is a CAKKT point.

Proof. For all k, we have $||y^k - x^k|| = ||y^k - (y^k + t_{k-1}d^{k-1})|| \le t_{k-1}||d^{k-1}||$. As $d^{k-1} \to 0$ by Lemma 5.7 and $t_{k-1} \le 1$, we obtain $||y^k - x^k|| \to 0$, and thus every limit point of $\{x^k\}$ is also a limit of $\{y^k\}$ and vice versa. So, it is sufficient to consider a limit point of $\{x^k\}$, let us say, $\bar{x} = \lim_{k \in K} x^k$.

By Lemma 5.5, $\phi(x^k) \to 0$, and thus \bar{x} is feasible. By Lemmas 5.4 and 5.7,

$$\|\operatorname{proj}_{\mathbb{L}^+_{\Omega}(y^k)}(0,0,-\nabla_x \mathcal{L}(y^k,\lambda^k,\mu^k))\| \le \sigma \|(s^{k-1},\ell_{k-1},d^{k-1})\| + \varepsilon_{k-1} \to 0.$$

Step 2 ensures that $\lim_{k \in K} \mu_j^k = 0$ whenever $g_j(\bar{x}) = \lim_{k \in K} g_j(y^k) < 0$. Thus, Corollary 3.2 implies that \bar{x} is a CAKKT point associated with the sequence $\{y^k\}_{k \in K}$.

Given a feasible point x that satisfies an optimality sequential condition, one may ask what the weakest property for x is to be a KKT point, for every objective function that has x as a local minimizer. Such a property is known as weakest strict CQ (it is indeed a CQ). In other words, the weakest strict CQs play the role of Guignard's CQ for sequential optimality conditions. In [5], the weakest strict CQs for several sequential conditions were provided. In particular, the weakest strict CQ associated with CAKKT is known as CAKKT-regularity, which we recall in the sequel. For $x \in \mathbb{R}^n$ and $\alpha \geq 0$, define the set-valued mapping by

$$K_{\bar{x}}(x,\alpha) = \left\{ \nabla h(x)\lambda + \nabla g(x)\mu \, \middle| \, \begin{array}{l} \displaystyle \sum_{i=1}^m |h_i(x)\lambda_i| + \sum_{j=1}^p |g_j(x)\mu_j| \leq \alpha, \\ \mu \geq 0, \quad \mu_j = 0, \ \forall j \not\in A(\bar{x}) \end{array} \right\}.$$

Definition 5.1. We say that a feasible \bar{x} for (NLP) satisfies the CAKKT-regular condition if $\limsup_{(x,\alpha)\to(\bar{x},0)} K_{\bar{x}}(x,\alpha) \subset K_{\bar{x}}(\bar{x},0)$.

CAKKT-regular is indeed a CQ since it implies Abadie's CQ [5, Theorem 6]. On the other hand, it is implied by very mild CQs in the literature (for a complete relationship between various CQs, see [5, Figure 6]). The next result is a direct consequence of Theorem 5.8 and [5, Theorem 2].

Corollary 5.9. Under the assumptions of Theorem 5.8, every limit point generated by Algorithm 1 that satisfies the CAKKT-regular condition is a KKT point.

At this point, we have shown that the new optimization step (Step 3 of Algorithm 1) allows us to ensure convergence to CAKKT points, improving the global convergence results of previous methods. The attentive reader may ask whether the new optimization step is really necessary to achieve CAKKT points in the IR framework. The answer is yes. We illustrate this by means of a simple example that our linearization of the feasible set (4), where the auxiliary variables s and ℓ are introduced, leads, under hypotheses H1 to H4, to an IR method with stronger convergence than the usual ones using linearization (3). In other words, keeping s and ℓ in the quadratic subproblems of Algorithm 1 is theoretically better than taking $(s,\ell)=0$.

Example 5.1. Let us consider the example [6]

min
$$\frac{(x_2-2)^2}{2}$$
 s.t. $x_1=0$, $x_1x_2=0$.

The unique minimizer is (0,2). Hypothesis H1 is immediate and H2 is fulfilled, for example for $H_k = I$. Taking $\lambda_{\min} = \lambda_{\max} = \mu_{\max} = 0$, (14) takes the form

min
$$(y_2^{k+1} - 2)d_2 + \frac{1}{2}d^T I d + \frac{1}{2}(s_1^2 + s_2^2)$$

s.t. $s_1 y_1^{k+1} + d_1 = 0$, $s_2 y_1^{k+1} y_2^{k+1} + y_2^{k+1} d_1 + y_1^{k+1} d_2 = 0$,

from which we conclude that its solution (s^k, d^k) satisfies

$$d_1^k = -s_1^k y_1^{k+1}$$
 and $y_1^{k+1} d_2^k = (s_1 - s_2) y_1^{k+1} y_2^{k+1}$.

Let us analyze the case where s is not present, or equivalently, s = (0,0). From the above relations, $d^k = (0,0)$ whenever $y_1^{k+1} \neq 0$, and thus $x^{k+1} = y^{k+1}$. So, in this case we go back to the restoration phase, obtaining a new restoration point y^{k+2} such that $\phi(y^{k+2}) \leq r\phi(x^{k+1}) = r\phi(y^{k+1})$ with r = 1/2. Note that the value of θ_k (Step 4 of Algorithm 1) does not matter. Thus, the classical IR approach, without s, can generate the sequence $\{x^k = y^k = (1/2^k, 1)\}$ since

$$\phi(y^{k+2}) = \left\| \left(\frac{1}{2^{k+2}}, \frac{1}{2^{k+2}} \right) \right\| \leq \frac{1}{2} \left\| \left(\frac{1}{2^{k+1}}, \frac{1}{2^{k+1}} \right) \right\| = \frac{1}{2} \phi(x^{k+1}), \quad \forall k \in \mathbb{N}.$$

Also, H3 and H4 are satisfied since $f(y^{k+1}) - f(x^k) = 0$ for all k. However, the limit point (0,1) is not the solution, or even a CAKKT point [6].

On the other hand, Algorithm 1, with possibly nonzero s, can not converge to (0,1) with a sequence satisfying H3 and H4 due to Theorem 5.8.

We end this section with an interesting consequence of our theory. Note that if the point y^{k+1} from Step 1 of Algorithm 1 is feasible, then trivially $(s^k, \ell_k) = 0$ at the optimality of problem (14), since in this case we have $h(y^{k+1}) = 0$ and $g^+(y^{k+1}) = 0$. Thus, we can remove the variables s and ℓ from (14). So, under hypotheses H1–H4, Theorem 5.8 says that the usual IR (Algorithm 1 without the variables s and ℓ) converges to CAKKT points whenever only feasible points y^{k+1} are computed. We will refer to such method as exact restoration. We summarize this result below.

Corollary 5.10. Suppose valid H1–H4 and consider the infinite sequences $\{x^k\}$, $\{y^k\}$ generated by the exact restoration method, i.e., Algorithm 1 without s and ℓ , and where y^{k+1} is always feasible. Then every limit point of $\{x^k\}$ or $\{y^k\}$ is CAKKT. Furthermore, if such limits satisfy the CAKKT-regular condition, then they are KKT.

Note that in Example 5.1, the sequence $\{y^{k+1}\}$ is infeasible. Exact restoration can only be expected on very specific problems where feasibility is easy to achieve (e.g., by closed formulas). The hard-spheres problem is an example [30]. We emphasize, however, that achieving exact feasibility is quite unusual for general problems. Our approach, in turn, recovers good theoretical convergence with a general inexact restoration phase.

6 Numerical tests

We implemented Algorithm 1 in Julia (v1.8.5) [10]. In this section, we discuss how the steps of the method were addressed in practice. Numerical tests on CUTEst problems are reported.

6.1 Practical aspects and implementation issues

6.1.1 Restoration step

The point y^{k+1} must satisfy $\phi(y^{k+1}) \leq r\phi(x^k)$. If $\phi(x^k)$ is small enough $(\phi(x^k) \leq \varepsilon_{\text{feas}})$, we choose $y^{k+1} = x^k$. Otherwise, we apply the Barzilai-Borwein-like method ABBmin₁ [23] with projections on box constraints, as done in [15], to the infeasibility problem

$$\min_{y,w} \frac{1}{2} \|h(y)\|^2 + \frac{1}{2} \|g(y) - w\|^2 + \frac{\xi}{2} \|y - x^k\|^2 \quad \text{s.t.} \quad w \le 0,$$
(28)

 $\xi \geq 0$, until $\phi(y^{k+1}) \leq r\phi(x^k)$ is reached or the projected gradient sup-norm (of its objective function) if less than or equal to 10^{-14} . We initialize $\xi = 10^{-4}$ and decrease it $(\xi \leftarrow \xi/10)$ whenever the projected gradient sup-norm is less than or equal to 10^{-7} , $\phi(y) > r\phi(x^k)$ and $\xi > 10^{-16}$. If ABBmin₁ fails, we proceed with the best-found point. Finally, if the infeasibility does not improve during ten consecutive outer iterations, we stop Algorithm 1 declaring failure. This criterion, based solely on numerical practice, is more flexible than the criterion in Step 1.

6.1.2 Estimation of Lagrange multipliers

We initialize $(\lambda^0, \mu^0) = (0, 0)$. New multipliers $(\lambda^{k+1}, \mu^{k+1})$ can be obtained from the dual solution $(\lambda_{\text{QP}}^{k-1}, \mu_{\text{QP}}^{k-1})$ of the quadratic problem of the previous iteration k-1. In fact, from the optimality conditions of (14), we have

$$\nabla_x \mathcal{L}(z^k) + H_{k-1} d^{k-1} + \sum_i \lambda_{\text{QP},i}^{k-1} \nabla h_i(y^k) + \sum_j \mu_{\text{QP},j}^{k-1} \nabla g_j(y^k) = 0, \quad \mu_{\text{QP}}^{k-1} \ge 0,$$

where $z^k := (y^k, \lambda^k, \mu^k)$. So, when $||H_{k-1}d^{k-1}||_{\infty}$ is small the safeguarded multiplier vector λ^{k+1} of iteration k can be approximated by

$$\lambda_i^{k+1} = \max\{\lambda_{\min}, \min\{\lambda_{\max}, \lambda_i^k + \lambda_{\mathrm{QP},i}^{k-1}\}\},\tag{29}$$

for all i (analogously for μ_j^{k+1}). The use of multiplier estimates from the subproblem was proposed, for example, in [16]. Unfortunately, the estimative (29) may be poor if $||H_{k-1}d^{k-1}||$ is large. We then use it if it promotes a reduction of $||\nabla_x \mathcal{L}(y^{k+1}, \cdot, \cdot)||_{\infty}$, otherwise we keep the multipliers unchanged. Although in the early stages of the optimization process the quantity $||H_{k-1}d^{k-1}||$ tends to be large, our numerical experience indicates that (29) is effective in reducing $||\nabla_x \mathcal{L}||$. Thus, this is the main strategy for updating multiplier estimates.

The new multiplier μ^{k+1} may not satisfy the complementarity condition. If this is the case, we try to correct μ^{k+1} when the feasibility tolerance is reached at y^{k+1} (that is, $\min\{\phi(y^{k+1}), \|(h(y^{k+1}), g^+(y^{k+1}))\|_{\infty}\} \leq \varepsilon_{\text{feas}}$) but the complementarity is not $(\|\min\{-g(y^{k+1}), \mu^{k+1}\}\|_{\infty} > \varepsilon_{\text{compl}})$, by applying the strategy described next. Let us consider the least squares problem

$$\min_{\lambda,\mu} \left\| \begin{bmatrix} \nabla h(y^{k+1}) & \nabla g(y^{k+1}) \\ 0 & \rho G(y^{k+1}) \end{bmatrix} \begin{bmatrix} \lambda \\ \mu \end{bmatrix} + \begin{bmatrix} \nabla f(y^{k+1}) \\ 0 \end{bmatrix} \right\|^2 \text{ s.t. } \mu \ge 0, \tag{30}$$

where $\rho > 0$ and $G(y^{k+1})$ is the diagonal matrix formed from $g^-(y^{k+1})$. The scalar ρ can be viewed as a penalization parameter for the products $\mu_j g_j^-(y^{k+1})$. As ρ grows, complementarity tends to be satisfied; on the other hand, optimality tends to be lost. We then adopt the following two-phase procedure:

- 1. we solve the unconstrained version of (30) starting with $\rho = 0.1$ using a suitable factorization. Then we check if the complementarity measure of Step 2 is at most $\max\{\varepsilon_{\text{compl}}, 0.1 \cdot \|\min\{-g(y^{k+1}), \mu^{k+1}\}\|_{\infty}\}$. If yes, we stop declaring "success" if the new μ is non-negative, otherwise we declare "failure". In the case of complementarity was not reduced enough, we check if the optimality measure (the sup-norm of the first row in (30)) is greater than $0.999 \cdot \|\nabla_x \mathcal{L}(y^{k+1}, \lambda^{k+1}, \mu^{k+1})\|_{\infty}$. If yes, we stop declaring "failure" and return the current multiplier estimate. Otherwise, we update $\rho \leftarrow 10\rho$ and repeat the process until complementarity was reduced or $\rho \geq 10^{20}$;
- 2. if the first phase fails, we simply reset all the multipliers to zero.

Remark 4. When (NLP) has only equality constraints, (30) is an unconstrained problem. In this case, no complementarity should be verified. However, it was observed numerically that even in this case recomputing multipliers can be useful to improve optimality. Therefore, we apply the above strategy if also the feasibility tolerance is reached and $\|\nabla_x \mathcal{L}(y^{k+1}, \lambda^{k+1}, \mu^{k+1})\|_{\infty} > \varepsilon_{\text{opt}}$. Note that in equality-constrained problems, the solution of (30) is always accepted as new the multiplier vector λ^{k+1} .

6.1.3 Optimization step

To simplify our exposition, we start by assuming that (NLP) has only equality constraints and that $H_k = \nabla_{xx}^2 \mathcal{L} + \sigma I_n$ for some $\sigma \geq 0$. The KKT system of the corresponding quadratic problem from step 3 is (we omit y^{k+1} , λ^{k+1} for clarity)

$$\begin{bmatrix} \nabla_{xx}^{2} \mathcal{L} + \sigma I_{n} & 0 & \nabla h \\ 0 & I_{m} & \operatorname{diag}(h) \\ \hline \nabla h^{T} & \operatorname{diag}(h) & -\xi I_{m} \end{bmatrix} \begin{bmatrix} d \\ s \\ \hline \lambda_{\mathrm{QP}} \end{bmatrix} = \begin{bmatrix} -\nabla_{x} \mathcal{L} \\ 0 \\ \hline 0 \end{bmatrix}$$
(31)

with $\xi=0$, where diag(h) is the diagonal matrix formed from $h(y^{k+1})$. Let M be the 2×2 block formed from $\nabla^2_{xx}\mathcal{L} + \sigma I_n$ and I_m . Considering the hypothesis H2', we want to compute $\sigma\geq 0$ such that M is positive definite on the kernel of $[\nabla h^T \operatorname{diag}(h)]$ (note that it is not necessary to adjust the identity I_m in M). A practical way to do this is increasing successively σ until the inertia of the matrix of coefficients of (31) is correct, that is, when such matrix has n+m positive and m negative eigenvalues. Also, it is necessary to correct a possible deficient rank of $[\nabla h^T \operatorname{diag}(h)]$, which can be done taking a positive ξ , see [33]. Similar to [11] (see also [33]), we compute the inertia of the matrix of coefficients of (31) by an inertia-revealing factorization. We adopt the following procedure: we start trying $\sigma=0$ and $\xi=0$. If the inertia is correct, we compute the solution using the available factorization. Otherwise, we update $\xi \leftarrow \max\{10^{-8}, 3\xi\}$ if the number of negative eigenvalues is less than m and $\sigma \leftarrow \max\{10^{-8}, 3\sigma\}$ if the number of positive eigenvalues is less than n+m, and recompute the factorization. This procedure is repeated until the inertia is correct. It is worth mentioning that, although H2' involves the lower bound σ_{\min} on the positivity of H_k , we allow H_k being $\nabla^2_{xx}\mathcal{L}$ without further verification in order to favour Newton-type steps.

For the case of equality and inequality constraints, we implemented a simple primal activeset method, see [33]. At each face, we need to solve an equality-constrained quadratic problem compose by all equalities and by the inequality constraints with indices in a working set W written as equalities, which leads to the coefficient matrix

$$\begin{bmatrix} \nabla_{xx}^{2} \mathcal{L} + \sigma I_{n} & 0 & 0 & \nabla h & \nabla g_{W} \\ 0 & I_{m} & 0 & \operatorname{diag}(h) & 0 \\ 0 & 0 & 1 & 0 & (g_{W}^{+})^{T} \\ \hline \nabla h^{T} & \operatorname{diag}(h) & 0 & -\xi I_{m} & 0 \\ \nabla g_{W}^{T} & 0 & g_{W}^{+} & 0 & -\xi \end{bmatrix}$$
(32)

where g_W^+ is the column vector formed by $g_j^+(y^{k+1})$, $j \in W$. Each quadratic problem is solved analogously to the equality-constrained case. We start with $W = \emptyset$ and compute $\sigma \geq 0$ as before. This σ serves for all subsequent subproblems, as equality constraints are present in all of them.

Remark 5. When $\nabla_x \mathcal{L}(y^{k+1}, \lambda^{k+1}, \mu^{k+1}) = 0$, (14) does not need to be solved since trivially $(s^k, \ell_k, d^k) = 0$ is its solution. This is in accordance with the purpose of d^k , which is intended to be a descent direction for the Lagrangian function at y^{k+1} . So, we pass directly to Step 1 whenever $\|\nabla_x \mathcal{L}(y^{k+1}, \lambda^{k+1}, \mu^{k+1})\|_{\infty} \leq 10^{-16}$.

Our implementation handles bounds on variables by treating them as ordinary inequality constraints, except in the restoration phase, where they are handled directly by introducing box-constraints $l \leq y \leq u$ in the infeasibility problem (28). In general, solving (14) deserves further research; for example, dual active-set methods may be preferable when the number of constraints is large.

6.1.4 Globalization step

We implemented the simple backtracking procedure $t_k \leftarrow 0.5t_k$, starting with $t_k \leftarrow 1$, to achieve (15) and (16). Under the hypotheses of Lemma 5.2, this procedure is finite (for safety, we stop it when $y^{k+1} + t_k d^k$ is numerically equal to y^{k+1}). Note that if $t_k < 1$ is the accepted step-size and \hat{t} is the rejected step-size immediately before, then $\hat{t} \in [t_k, 2t_k]$.

6.2 Tests with problems from CUTEst

Tests were run on a computer equipped with Intel© i5-4570 CPU 3.20GHz, GNU/Linux Ubuntu 20.04.2 LTS. We set r=0.2, $\theta_{-1}=0.9999$, $\lambda_{\min}=-10^{20}$, $\mu_{\max}=\lambda_{\max}=10^{20}$, $\gamma_c=10^{-4}$, $\sigma_{\min}=10^{-8}$ and the maximum number of outer iterations equals 100. Supported by Corollary 3.2, the stopping criterion is

$$\max\{\|(s^k, \ell_k, d^k)\|_{\infty}\} \le \varepsilon_{\text{opt}}, \ \min\{\phi(y^{k+1}), \|(h(y^{k+1}), g^+(y^{k+1}))\|_{\infty}\} \le \varepsilon_{\text{feas}}$$
 (33)

where $\varepsilon_{\rm opt} = 10^{-6}$ and $\varepsilon_{\rm feas} = 10^{-8}$ (note that from Step 3 of Algorithm 1, (s^k, ℓ_k, d^k) is associated with y^{k+1}). In the restoration phase, we set the maximum number of ABBmin₁ iterations to 50,000. Also, we stop it due to lack of progress if no improvement in the objective function of (28) occurs during 300 consecutive iterations.

We consider the following two IR strategies:

- IR_{0,0}: the usual IR, where we eliminate the auxiliary variables s and ℓ from the quadratic subproblem (14) (so, it is the same as fixing $(s,\ell) = 0$ in Algorithm 1). This corresponds to projecting onto $\mathbb{L}_{\Omega}(y^{k+1})$, see (3);
- $IR_{s,\ell}$: Algorithm 1 where (14) is solved considering s and ℓ .

We consider CUTEst problems with the number of constraints between 1 and 10,000. We excluded feasibility problems, i.e., those with constant objective function. This is because we initialize all multipliers to zero, so we always have $(\lambda^k, \mu^k) = 0$ and $\nabla_x \mathcal{L}(y^k, \lambda^k, \mu^k) = 0$. Therefore, (14) does not need to be solved (see Remark 5), i.e., $IR_{0,0}$ and $IR_{s,\ell}$ behave exactly the same.

6.2.1 Measuring the cost of the new optimization step

We are interested in determining whether the new optimization phase adds relevant computational effort compared to the usual one, that is, we want to know if solving (14) is significantly more expensive than solving the same quadratic problem without s and ℓ . Firstly, it is worth noting that compared to the quadratic problem without the variable s, the matrix (31) aggregates only the two diagonal blocks I_m and diag(h). This results in an additional storage requirement of only 2m elements. The same applies when inequality constraints are present.

As the restoration phase represents an important portion of the total execution time that can vary depending on the generated sequence $\{x^k\}$, we performed only one iteration of ABBmin₁ for a fairer comparison (of course, we are not analyzing robustness at this time). We also disabled the computation of multiplier estimate (section 6.1.2) by imposing $\mu_{\text{max}} = \lambda_{\text{min}} = 0$ and we limited the number of outer iterations to 5. To measure the effort on a more representative set of problems, we select only those with 1,000 to 10,000 equality constraints and no inequalities (when only equality constraints are present, the optimization phase boils down to solving exactly one system (31)). In the 39 selected problems, the runtime of $IR_{s,\ell}$ and $IR_{0,0}$ were essentially the same, that is, the costs of solving the optimality system with matrices (31) and (32) are essentially the same. Table 1 shows the runtime comparison for the selected problems grouped by the number of equality constraints (m) – note that for each equality constraint, there is one additional variable s_i in $IR_{s,\ell}$. The first two columns indicate the interval of m considered and the quantity of problems within each of them. The third column is obtained as follows: for each problem p, we define the runtime $T_{s,\ell}^p$ of $IR_{s,\ell}$ as the average of the execution times of the necessary runs to reach 10 seconds; $T_{0,0}^p$ is defined analogously. Then, we compute the geometric mean of $T_{s,\ell}^p/T_{0,0}^p$ across all problems p in the interval. This serves to measure the overall percentage of $IR_{s,\ell}$'s runtime relative to $IR_{0,0}$, which is reported in the fourth column. The last column is the standard deviation of $\max\{0, T_{s,\ell}^p - T_{0,0}^p\}$ across all problems in the interval; it aims to measure the dispersion in runtime between the problems where $IR_{s,\ell}$ was slower than $IR_{0,0}$. From our tests, we can say that the addition of variable s in (31) does not affect the cost of computing a direction in the optimization step (we attribute the minor fluctuations in the data presented in Table 1 to natural variations in execution time). Also, this conclusion remains valid when the number of equality constraints varies. Although the number of variables increases along with the number of equality constraints, Table 2 brings the problems grouped by the number of variables, from which similar conclusion can be verified.

In the presence of inequality constraints we observed discrepancies in runtime either in favor of $IR_{s,\ell}$ or in favor of $IR_{0,0}$, and thus no consistent conclusion could be drawn. Although an arbitrary number of inequality constraints entails the addition of a single variable ℓ , which results in the larger matrix (32), the effect on the number of faces explored in the active-set method may be different from an IR variant to another. As detailed in the next section, handling inequality constraints in inexact restoration methods remains a challenge. It is worth mentioning that the implementation in [11] only handles equality-constrained problems with all free variables.

Table 1: Runtime comparison between problems only with equality constraints, categorized by the number of constraints. $IR_{0,0}$ is the reference.

number of	number of	geometric mean	$IR_{s,\ell}$ relative	standard deviation of
constraints	problems	of $T_{s,\ell}^{p}/T_{0,0}^{p}$	to $IR_{0,0}$	$\max\{0, T_{s,\ell}^p - T_{0,0}^p\}$
100-1000	4	1.0098	0.98% slower	0.04492
1001-2000	8	0.9998	0.02% faster	0.01250
2001-3000	4	0.9955	0.45% faster	0.00375
3001-4000	4	0.9998	0.02% faster	0.00419
4001-5000	5	1.0024	0.24% slower	0.00405
5001-7000	3	1.0009	0.09% slower	0.00231
7001-8000	7	0.9895	1.05% faster	0.03216
> 8000	4	0.9908	0.92% faster	0.01758

Table 2: Runtime comparison between problems only with equality constraints, categorized by the number of variables. $IR_{0,0}$ is the reference.

number of	number of	geometric mean	$IR_{s,\ell}$ relative	standard deviation of
variables	problems	of $T_{s,\ell}^{p}/T_{0,0}^{p}$	to $IR_{0,0}$	$\max\{0, T_{s,\ell}^p - T_{0,0}^p\}$
100-2000	5	1.0103	1.03% slower	0.09322
2001-3000	7	0.9981	0.19% faster	0.02132
3001-6000	4	0.9955	0.45% faster	0.00405
6001-9000	5	1.0040	0.40% slower	0.01059
9001-10000	13	0.9925	0.75% faster	0.02422
> 10000	5	0.9961	0.39% faster	0.02280

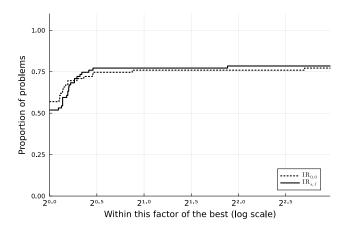


Figure 2: Comparison of the number outer iterations between $IR_{s,\ell}$ and $IR_{0,0}$ in equality-constrained problems.

6.2.2 Comparison of IR methods

Since the cost of solving one quadratic problem in Step 3 of Algorithm 1 is essentially the same for $IR_{s,\ell}$ and $IR_{0,0}$, to compare the IR methods we disregard possible variations in the performance of the solver employed in the restoration phase. Therefore, we compare $IR_{s,\ell}$ and $IR_{0,0}$ by the number of outer iterations.

We first present the results on problems with only equality constraints. From the 90 selected problems, 11 were discarded due to unexpected errors or 1 hour runtime exceeded. In general, both IR methods work similarly. From the 79 problems considered, IR_{0,0} and IR_{s,ℓ} solved 77.21% and 78.48%, respectively. For problems MSS2 and MSS3, IR_{0,0} and IR_{s,ℓ} stopped at a stationary point of the infeasibility, respectively. We consider a point to be stationary for infeasibility if it fails to meet the feasibility criterion in (33), but is stationary for (28) with $\xi = 0$. The same situation occurred for problems MSS1 and S308NE with both strategies. IR_{s,ℓ} needs more iterations to converge than IR_{0,0} in 42.37% of the problems that were solved by both methods. However, in 36.00% of them only one additional iteration was performed; in 28.00% two iterations, in 24.00% between three and six iterations and in 12.00%, 35 or more iterations. On the other hand, among problems where IR_{s,ℓ} converged with fewer iterations than IR_{0,0} (32.20% of the solved problems), we observed larger discrepancies favourable to IR_{s,ℓ}: HS56 required 66 fewer iterations and LUKVLE13, 15; LUKVLE12 and SPIN20P were solved by IR_{s,ℓ}, but not by IR_{0,0}. Figure 2 shows the performance profile [19] on problems with equality constraints only.

When inequality constraints are present, a difficulty in implementing robust inexact restoration algorithms lies in updating the multipliers in Step 2 due to complementarity. The strategy described in Section 6.1.2 attempts to overcome this issue, but it is certainly not optimal. In particular, when complementarity is not achieved, we reset all multipliers to zero, which, although theoretically valid, leads to a very poor numerical performance. In fact, the lack of accurate estimates for the Lagrange multipliers is the main reason for failure in our numerical tests: when $\lambda^{k+1} = 0$ and $\mu^{k+1} = 0$ in Step 2 of Algorithm 1, the direction computed in Step 3 tends to decrease only f, resulting in very small steps t_k close to the solution (even if $H_k \approx \nabla^2 f(y^{k+1})$, i.e., $\sigma_k \approx 0$), and in an optimality measure $\|\nabla \mathcal{L}(x^{k+1},0,0)\|_{\infty}$ far from zero. In this case, the method adopted in the feasibility phase is not capable of driving the minimization process to the solution by itself. From the 601 selected problems with inequality constraints and/or bounds on

variables (bounds are treated as inequality constraints in our implementation), the IR methods stopped within 1 hour in 420 of them. Among these problems, $IR_{0,0}$ converged in 53.57% of them, while $IR_{\ell,s}$ solved 54.05% of them. The relationship between outer iterations of the two variants is similar to the equality-constrained case (see Figure 2), that is, $IR_{\ell,s}$ typically requires a few more iterations than $IR_{0,0}$ to converge, but it is slightly more robust.

Considering that the cost of solving quadratic subproblems in both IR methods is similar and that the new IR method has shown a slight improvement in robustness, we believe that future implementations can benefit from using the new optimization step. In particular, the resolution of quadratic subproblems is a crucial issue (the active set method failed in 36 problems). Finally, we executed the IR methods on problems with inequality constraints reformulated using slacks variables; this is common in the literature [9, 11, 13, 16, 21, 29]. We observed a worse performance compared to treating inequalities directly. However, it is worth noting that this outcome can vary depending on the implementation, especially in terms of how the quadratic problems in the optimization step are solved.

Following a reviewer's suggestion, we conducted all tests with a limit of 1,000 outer iterations. The scenario for equality-constrained problems remained almost unchanged (only one more problem was solved by the new IR algorithm). For problems with inequality constraints and/or bounds on variables, both IR variants solved 10 additional problems, while 3 were solved only by the standard IR and another 4 only by the new IR. This represents only 2.8% of the total problems with inequalities. Thus, the gain from increasing the maximum number of iterations is marginal, and it does not alter the relative comparison between the two variants. The poor robustness exhibited by both IR methods on inequality-constrained problems suggests that specialized techniques for dealing with such constraints should be a topic of future research, as well as an effective strategy to compute the multiplier μ^{k+1} in Step 2. Some possible improvements are listed in the next section.

7 Conclusion

Sequential optimality conditions are powerful tools for unifying and establishing convergence of several numerical methods in optimization. In this work, we demonstrate how this concept can be used as a guide for the development of new inexact restoration algorithms with a focus on improving the global convergence theory. The fundamental idea behind our approach is to rewrite the CAKKT condition [6] as a projection onto a suitable linearization of the feasible set and incorporate it into a novel optimization step. This new step not only involves an optimality measure, but also carries additional information regarding the complementarity, a crucial ingredient to reach CAKKT points. This is accomplished by considering the level of infeasibility at the target point through auxiliary variables. To the best our knowledge, the resulting IR method is the first numerical method for standard nonlinear optimization that converges to CAKKT points without any exogenous assumptions while allowing inexact resolution of subproblems. Furthermore, our theory enables us to establish a new CAKKT convergence status for previous IR algorithms when exact feasibility is maintained throughout the optimization process. This finding helps explain why inexact restoration techniques work well in problems where feasibility can be achieved with high precision, for example the hard-spheres problem [27, 30]. However, despite all the theoretical apparatus, the computational performance in the numerical experiments of both the traditional IR and the new proposal did not present good robustness for general-purpose problems from CUTEst.

We compared our new IR strategy to that one with the classical, well-established, optimization step on CUTEst problems, concluding that the use of the new optimization step does not

imply a significant increase in runtime. This encourages us to pursue a state-of-art implementation. Based on our numerical experience, the implementation of Algorithm 1 and other IR methods can be improved considering, among others, the following topics: (i) adopting our novel optimization step as standard; (ii) employing a variety of specialized methods to compute points in Step 1 for different types of constraints; (iii) refining estimates of Lagrange multipliers, especially those associated with inequality constraints; (iv) the use of regularized quadratic subproblems as done in SQP methods [26]; (v) incorporating extrapolation steps and non-monotone line searches in Step 5 to avoid excessively small steps; (vi) the possibility of using a filter approach [27]; and (vii) dealing with variable bounds directly within the quadratic subproblems. Certainly, this involves a much careful implementation that should be considered in a future research, and some of these improvements require adjustments to the theory. Special attention should be devoted to specific problems where the inexact restoration philosophy has favourable properties, such as bilevel optimization [1, 7].

References

- [1] R. Andreani, S. L. C. Castro, J. L. Chela, A. Friedlander, and S. A. Santos. An inexact-restoration method for nonlinear bilevel programming problems. *Computational Optimization and Applications*, 43:307–328, 2009.
- [2] R. Andreani, G. Haeser, and J. M. Martínez. On sequential optimality conditions for smooth constrained optimization. *Optimization*, 60(5):627–641, 2011.
- [3] R. Andreani, G. Haeser, L. M. Mito, A. Ramos, and L. D. Secchin. On the best achievable quality of limit points of augmented Lagrangian schemes. *Numerical Algorithms*, 90:851– 877, 2022.
- [4] R. Andreani, G. Haeser, M. L. Schuverdt, and P. J. S. Silva. A relaxed constant positive linear dependence constraint qualification and applications. *Mathematical Programming*, 135(1):255–273, 2012.
- [5] R. Andreani, J. M. Martínez, A. Ramos, and P. J. S. Silva. Strict constraint qualifications and sequential optimality conditions for constrained optimization. *Mathematics of Operations Research*, 43(3):693–717, 2018.
- [6] R. Andreani, J. M. Martínez, and B. F. Svaiter. A new sequential optimality condition for constrained optimization and algorithmic consequences. SIAM Journal on Optimization, 20(6):3533–3554, 2010.
- [7] R. Andreani, V. A. Ramirez, S. A. Santos, and L. D. Secchin. Bilevel optimization with a multiobjective problem in the lower level. *Numerical Algorithms*, 81(3):915–946, 2019.
- [8] R. Andreani, A. Ramos, A. A. Ribeiro, L. D. Secchin, and A. R. Velazco. On the convergence of augmented Lagrangian strategies for nonlinear programming. *IMA Journal of Numerical Analysis*, 42(2):1735–1765, 2022.
- [9] M. B. Arouxét, N. E. Echebest, and E. A. Pilotta. Inexact restoration method for nonlinear optimization without derivatives. *Journal of Computational and Applied Mathematics*, 290:26–43, 2015.

- [10] J. Bezanson, A. Edelman, S. Karpinski, and V. B. Shah. Julia: A fresh approach to numerical computing. *SIAM Review*, 59(1):65–98, 2017.
- [11] E. Birgin, L. Bueno, and J. Martínez. Assessing the reliability of general-purpose inexact restoration methods. *Journal of Computational and Applied Mathematics*, 282:1–16, 2015.
- [12] E. G. Birgin, N. Krejic, and J. M. Martínez. Iteration and evaluation complexity for the minimization of functions whose computation is intrinsically inexact. *Mathematics of Computations*, 89:253–278, 2020.
- [13] E. G. Birgin and J. M. Martínez. Local convergence of an inexact-restoration method and numerical experiments. *Journal of Optimization Theory and Applications*, 127(2):229–247, 2005.
- [14] E. G. Birgin and J. M. Martínez. Practical Augmented Lagrangian Methods for Constrained Optimization. Society for Industrial and Applied Mathematics, Philadelphia, PA, 2014.
- [15] E. G. Birgin, J. M. Martínez, and M. Raydan. Nonmonotone spectral projected gradient methods on convex sets. SIAM Journal on Optimization, 10(4):1196–1211, 2000.
- [16] L. F. Bueno, G. Haeser, and J. M. Martinez. A flexible inexact restoration methods for constrained optimization. *Journal on Optimization Theory and Applications*, 165:188–208, 2015.
- [17] L. F. Bueno, G. Haeser, and J. M. Martínez. An inexact restoration approach to optimization problems with multiobjective constraints under weighted-sum scalarization. *Optimization Letters*, 10:1315–1325, 2016.
- [18] L. F. Bueno and J. M. Martínez. On the complexity of an inexact restoration method for constrained optimization. SIAM Journal on Optimization, 30(1):80–101, 2020.
- [19] E. D. Dolan and J. J. Moré. Benchmarking optimization software with performance profiles. *Mathematical Programming*, 91(2):201–213, 2002.
- [20] A. Fischer and A. Friedlander. A new line search inexact restoration approach for nonlinear programming. *Computational Optimization and Applications*, 46(2):333–346, 2010.
- [21] J. B. Francisco, D. S. Gonçalves, F. S. V. Bazán, and L. L. T. Paredes. Non-monotone inexact restoration method for nonlinear programming. *Computational Optimization and Applications*, 76(3):867–888, 2019.
- [22] J. B. Francisco, J. M. Martínez, L. Martinez, and F. I. Pisnitchenko. Inexact restoration methods for minimization problems that arise in electronic structure calculations. *Computational Optimization and Applications*, 50:555–590, 2011.
- [23] G. Frassoldati, L. Zanni, and G. Zanghirati. New adaptive stepsize selections in gradient methods. *Journal of Industrial and Management Optimization*, 4:299–312, 2008.
- [24] P. E. Gill, V. Kungurtsev, and D. P. Robinson. A shifted primal-dual penalty-barrier method for nonlinear optimization. *SIAM Journal on Optimization*, 30(2):1067–1093, 2020.
- [25] C. C. Gonzaga, E. Karas, and M. Vanti. A globally convergent filter method for nonlinear programming. SIAM Journal on Optimization, 14(3):646–669, 2004.

- [26] A. F. Izmailov and M. V. Solodov. Stabilized SQP revisited. *Mathematical Programming*, 133(1):93–120, 2012.
- [27] E. W. Karas, E. A. Pilotta, and A. A. Ribeiro. Numerical comparison of merit function with filter criterion in inexact restoration algorithms using hard-spheres problems. *Computational Optimization and Applications*, 44(3):427–441, 2009.
- [28] C. Y. Kaya. Inexact restoration for Runge–Kutta discretization of optimal control problems. SIAM Journal on Numerical Analysis, 48(4):1492–1517, 2010.
- [29] J. M. Martínez. Inexact-restoration method with Lagrangian tangent decrease and new merit function for nonlinear programming. *Journal of Optimization Theory and Applica*tions, 111(1):39–58, 2001.
- [30] J. M. Martínez and E. A. Pilotta. Inexact-restoration algorithm for constrained optimization. *Journal of Optimization Theory and Applications*, 104(1):135–163, 2000.
- [31] J. M. Martinez and E. A. Pilotta. Inexact restoration methods for nonlinear programming: Advances and perspectives. In L. Qi, K. Teo, and X. Yang, editors, *Optimization and Control with Applications*, pages 271–291, Boston, MA, 2005. Springer US.
- [32] J. M. Martínez and B. F. Svaiter. A practical optimality condition without constraint qualifications for nonlinear programming. *Journal of Optimization Theory and Applications*, 118(1):117–133, 2003.
- [33] J. Nocedal and S. J. Wright. Numerical Optimization. Springer, New York, 2 edition, 2006.
- [34] R. T. Rockafellar and R. J.-B. Wets. *Variational Analysis*, volume 317 of *Grundlehren der mathematischen Wissenschaften*. Springer-Verlag Berlin Heidelberg, 1 edition, 1998.