

Complementarity Formulations of ℓ_0 -norm Optimization Problems

Mingbin Feng^{*}, John E. Mitchell[†], Jong-Shi Pang[‡], Xin Shen[§], and Andreas Wächter[¶]

September 25, 2013

Abstract

In a number of application areas, it is desirable to obtain sparse solutions. Minimizing the number of nonzeroes of the solution (its ℓ_0 -norm) is a difficult nonconvex optimization problem, and is often approximated by the convex problem of minimizing the ℓ_1 -norm. In contrast, we consider exact formulations as mathematical programs with complementarity constraints and their reformulations as smooth nonlinear programs. We discuss properties of the various formulations and their connections to the original ℓ_0 -minimization problem in terms of stationarity conditions, as well as local and global optimality. Numerical experiments using randomly generated problems show that standard nonlinear programming solvers, applied to the smooth but nonconvex equivalent reformulations, are often able to find sparser solutions than those obtained by the convex ℓ_1 -approximation.

Keywords: ℓ_0 -norm minimization, complementarity constraints, nonlinear programming

1 Introduction

Denoted by $\|\bullet\|_0$, the so-called ℓ_0 -norm of a vector is the number of nonzero components of the vector. In recent years, there has been an increased interest in solving optimization problems that minimize or restrict the number of nonzero elements of the solution vector [1, 2, 5, 7, 8, 9, 27, 30]. A simple example of such a problem is that of finding a solution to a system of linear inequalities

^{*}Department of Industrial Engineering and Management Sciences, Northwestern University, Evanston, IL 60208, USA. This author is supported by National Science Foundation grant DMS-1216920. E-Mail: mingbinfeng2011@u.northwestern.edu

[†]Department of Mathematical Sciences, Rensselaer Polytechnic Institute, Troy, NY 12180, USA. This author was supported by the National Science Foundation under Grant Number CMMI-1334327 and by the Air Force Office of Scientific Research under Grant Number FA9550-11-1-0260. E-Mail: mitchj@rpi.edu

[‡]Department of Industrial and Systems Engineering, University of Southern California, Los Angeles, CA 90089, USA. This author was supported by the National Science Foundation under Grant Number CMMI-1333902 and by the Air Force Office of Scientific Research under Grant Number FA9550-11-1-0151. E-Mail: jongship@usc.edu

[§]Department of Mathematical Sciences, Rensselaer Polytechnic Institute, Troy, NY 12180, USA. This author was supported by the National Science Foundation under Grant Number CMMI-1334327 and by the Air Force Office of Scientific Research under Grant Number FA9550-11-1-0260. E-Mail: shenx5@rpi.edu

[¶]Department of Industrial Engineering and Management Sciences, Northwestern University, Evanston, IL 60208, USA. This author was supported by National Science Foundation grant DMS-1216920 and grant CMMI-1334639. Email: andreas.waechter@northwestern.edu

with the least ℓ_0 -norm:

$$\begin{aligned} & \underset{x \in \mathbb{R}^n}{\text{minimize}} && \|x\|_0 \\ & \text{subject to} && Ax \geq b \quad \text{and} \quad Cx = d, \end{aligned} \tag{1}$$

where $A \in \mathbb{R}^{m \times n}$, $C \in \mathbb{R}^{k \times n}$, $b \in \mathbb{R}^m$ and $d \in \mathbb{R}^k$ are given matrices and vectors, respectively. Since this problem is NP-hard, one popular solution approach replaces the nonconvex discontinuous ℓ_0 -norm in (1) by the convex continuous ℓ_1 -norm, leading to a linear program:

$$\begin{aligned} & \underset{x \in \mathbb{R}^n}{\text{minimize}} && \|x\|_1 \\ & \text{subject to} && Ax \geq b \quad \text{and} \quad Cx = d. \end{aligned} \tag{2}$$

Theoretical results are known that provide sufficient conditions under which an optimal solution to (2) is also optimal to (1) [4, 11, 17, 31]. Yet these results are of limited practical value as the conditions can not easily be verified or guaranteed for specific realizations of (1); thus in general, optimal solutions to (2) provide suboptimal solutions to (1).

It is our contention that, from a practical perspective, improved solutions to (1) can be obtained by reformulating the ℓ_0 -norm in terms of complementarity constraints [23]. This leads to a linear program with linear complementarity constraints (LPCC) which can be solved with specialized algorithms that do not depend on the feasibility and/or boundedness of the constraints [18, 19]. In the event that bounds are known on the solutions of the problem, the LPCC can be further reformulated as a mixed-integer linear program (MILP). However, the solution of this MILP is usually too time-consuming for large instances.

As an alternative to the MILP approach, the LPCC can be expressed directly as a smooth continuous nonlinear program (NLP). It is the main purpose of this research to examine the quality of solutions computed by standard NLP solvers applied to these smooth reformulations of the ℓ_0 -norm. There are two properties of the NLP reformulations that make them difficult to solve. First, the NLPs are highly nonconvex, and, consequently, the solutions returned by the NLP solvers depend strongly on the starting point, because the NLP methods are typically only able to find local minimizers or Karush-Kuhn-Tucker (KKT) points, instead of global minimizers. Secondly, the NLPs are not well-posed in the sense that they do not satisfy the assumptions that are made usually for the convergence analysis of standard NLP algorithms, such as constraint qualifications. Nevertheless, our numerical results show that these methods often generate high-quality solutions for the ℓ_0 -norm problem (1).

The remainder of this paper is organized as follows. In Section 1.1 we present two basic complementarity formulations for the ℓ_0 -norm. One of them leads to an LPCC formulation of the problem (1) which is reformulated as a smooth NLP using different approaches, including a new construction based on squared complementarities. The other complementarity formulation results in a nonlinear program with bilinear, disjunctive constraints. These formulations are generalized to the nonlinear case in Section 2 where we introduce an NLP model whose objective comprises a weighted combination of a smooth term and a discontinuous ℓ_0 -term. This model is sufficiently broad to encompass many optimization problems that include applications arising from compressive sensing [5, 9], basis pursuit [2, 8], LASSO regression [27, 30], image deblurring [1], and the least misclassification (as opposed to the well-known least-residual) support-vector machine problem with a soft margin; the latter problem was first introduced by Mangasarian [7, 24].

To give some theoretical background for the expected convergence behavior for (local) NLP solvers, connections between the KKT points of the smooth formulations of the complementarity problems

and the original ℓ_0 -problem are established in Section 3. Further insights are obtained in Section 4 by considering ε -relaxations of the smooth NLP formulations.

The practical performance of standard NLP codes for the solution of ℓ_0 -minimization problems is assessed in Section 5. We present numerical results for an extensive set of computational experiments that show that the solutions obtained by some NLP formulations of the ℓ_0 -minimization are significantly better than those obtained from the convex ℓ_1 -formulation, often close to the globally optimal objective value. Conclusions and an outlook for future research are given in the final section.

1.1 Equivalent formulations

We start by introducing two basic ways to formulate the ℓ_0 -norm using complementarity constraints.

Full complementarity. A straightforward way of formulating the ℓ_0 -norm using complementarity constraints is to first express $x = x^+ - x^-$ with x^\pm being the non-negative and non-positive parts of x , respectively; this is followed by the introduction of a vector $\xi \in [0, 1]^n$ that is complementary to $|x|$, the absolute-value vector of x . This maneuver leads to the following formulation:

$$\begin{aligned} \underset{x, x^\pm, \xi}{\text{minimize}} \quad & \mathbf{1}_n^T (\mathbf{1}_n - \xi) = \sum_{j=1}^n (1 - \xi_j) \\ \text{subject to} \quad & Ax \geq b, \quad Cx = d, \quad \text{and} \quad x = x^+ - x^- \\ & 0 \leq \xi \perp x^+ + x^- \geq 0 \\ & 0 \leq x^+ \perp x^- \geq 0 \\ \text{and} \quad & \xi \leq \mathbf{1}_n, \end{aligned} \tag{3}$$

where $\mathbf{1}_n$ is the n -vector of all ones. It is not difficult to deduce that if x is an optimal solution of (1), then by letting $x^\pm \triangleq \max(0, \pm x)$ and

$$\xi_j \triangleq \begin{cases} 0 & \text{if } x_j \neq 0 \\ 1 & \text{if } x_j = 0 \end{cases} \quad j = 1, \dots, n, \tag{4}$$

the resulting triple (x^\pm, ξ) is an optimal solution of (3) with objective value equal to $\|x\|_0$. Conversely, if (x^\pm, ξ) is an optimal solution of (3), then $x \triangleq x^+ - x^-$ is an optimal solution of (1) with the same objective value as the optimal objective value of (3). The definition (4) provides a central connection between (1) and its “pieces” to be made precise in Section 3. Such pieces are smooth programs in which some of the x -variables are fixed at zero and correspond in some way to the enumeration of the zero versus nonzero components of x . The scalar $1 - \xi_i$ is the indicator of the support of x_i ; we call $\mathbf{1}_n - \xi$ the *support vector* of x .

It is easy to see that the complementarity between the variables x^\pm is not needed in (3); this results in the following equivalent formulation of this problem, and thus of (1):

$$\begin{aligned} \underset{x, x^\pm, \xi}{\text{minimize}} \quad & \mathbf{1}_n^T (\mathbf{1}_n - \xi) = \sum_{j=1}^n (1 - \xi_j) \\ \text{subject to} \quad & Ax \geq b, \quad Cx = d, \quad \text{and} \quad x = x^+ - x^- \\ & 0 \leq \xi \perp x^+ + x^- \geq 0 \\ & x^\pm \geq 0 \quad \text{and} \quad \xi \leq \mathbf{1}_n, \end{aligned} \tag{5}$$

Nevertheless, in terms of the global resolution of (3), maintaining the complementarity between x^\pm could potentially allow sharper cutting planes to be derived in a branch-and-cut scheme for solving this disjunctive program. This also led to better numerical results in our experiments reported in Section 5.

There are several well-known equivalent formulations of the complementarity condition $0 \leq y \perp z \geq 0$ in (3) and (5) that lead to a smooth continuous NLP formulation:

- $(y, z) \geq 0$ and $y^T z \leq 0$ (inner product complementarity);
- $(y, z) \geq 0$ and $y \circ z \leq 0$, where $u \circ v$ denotes the Hadamard, i.e., componentwise, product of two vectors u and v (componentwise or Hadamard complementarity);
- Adding the penalty term $My^T z$ in the objective (penalized complementarity) for some large scalar $M > 0$;
- $(y, z) \geq 0$ and $(y^T z)^2 \leq 0$ (squared complementarity).

Interestingly, the last formulation, which has never been used in the study of complementarity constraints, turns out to be quite effective for solving some ℓ_0 -norm minimization problems.

We point out that, with the exception of the penalizing complementarity approach, none of these reformulations of the complementarity problem give rise to a well-posed NLP model in the sense that the Mangasarian-Fromovitz constraint qualification (MFCQ) fails to hold at any feasible point, and the existence of KKT points is not guaranteed. Nevertheless, some NLP solvers have been found to be able to produce good numerical solutions for these reformulations [14].

Half complementarity. There is a simpler formulation, which we call the *half complementarity formulation*, that requires only the auxiliary ξ -variable:

$$\begin{aligned} & \underset{x, \xi}{\text{minimize}} && \mathbf{1}_n^T (\mathbf{1}_n - \xi) \\ & \text{subject to} && Ax \geq b; \quad Cx = d \\ & && 0 \leq \xi \leq \mathbf{1}_n; \quad \text{and} \quad \xi \circ x = 0, \end{aligned} \tag{6}$$

The equivalence of (1) and (6) follows from the same definition (4) of ξ . Strictly speaking, the constraints in (6) are not of the complementarity type because there is no non-negativity requirement on the variable x ; yet the Hadamard constraint $\xi \circ x = 0$ contains the disjunctions: either $\xi_i = 0$ or $x_i = 0$ for all i .

Finally, if a scalar $M > 0$ is known such that $M \geq \|x^*\|_\infty$ for an optimal solution x^* , then the ℓ_0 -norm minimization problem (1) can be formulated as a mixed-integer linear program with the introduction of a binary variable $\zeta \in \{0, 1\}^n$:

$$\begin{aligned} & \underset{x, \zeta}{\text{minimize}} && \mathbf{1}_n^T \zeta = \sum_{j=1}^n \zeta_j \\ & \text{subject to} && Ax \geq b \quad \text{and} \quad Cx = d, \\ & && -M\zeta \leq x \leq M\zeta, \\ & \text{and} && \zeta \in \{0, 1\}^n. \end{aligned} \tag{7}$$

2 A General ℓ_0 -norm Minimization Problem

Together, the ℓ_0 -norm and its complementarity formulation allow a host of minimization problems involving the count of variables to be cast as disjunctive programs with complementarity constraints. A general NLP model of this kind is as follows: for two finite index sets \mathcal{E} and \mathcal{I} ,

$$\begin{aligned} & \underset{x}{\text{minimize}} && f(x) + \gamma \|x\|_0 \\ & \text{subject to} && c_i(x) = 0, \quad i \in \mathcal{E} \\ & \text{and} && c_i(x) \leq 0, \quad i \in \mathcal{I}, \end{aligned} \tag{8}$$

where $\gamma > 0$ is a prescribed scalar and the objective function f and the constraint functions c_i are all continuously differentiable. A distinguished feature of this problem is that its objective function is discontinuous, in fact lower semicontinuous; as such, it attains its minimum over any compact set, but in general the existence/attainment of an optimal solution is not immediately clear. Among other things, the reformulation presented below offers a constructive venue for establishing the solvability of the problem, under reasonable conditions.

Similar to (3), we can derive the equivalent formulation as a complementarity problem:

$$\begin{aligned} & \underset{x, x^\pm, \xi}{\text{minimize}} && f(x) + \gamma^T (\mathbf{1}_n - \xi) \\ & \text{subject to} && c_i(x) = 0, \quad i \in \mathcal{E} \\ & && c_i(x) \leq 0, \quad i \in \mathcal{I} \\ & && x = x^+ - x^- \\ & && 0 \leq \xi \perp x^+ + x^- \geq 0 \\ & && 0 \leq x^+ \perp x^- \geq 0 \\ & \text{and} && \xi \leq \mathbf{1}_n, \end{aligned} \tag{9}$$

where we have used an arbitrary given positive vector γ instead of a scalar γ -multiple of the vector of ones. This time, the statement involves nonlinear objective and constraint functions, giving rise to a Mathematical Program with Complementarity Constraints (MPCC).

As for the half complementarity formulation (6) of (1), we may associate with (9) the following smooth NLP with an auxiliary variable ξ :

$$\begin{aligned} & \underset{x, \xi}{\text{minimize}} && f(x) + \gamma^T (\mathbf{1}_n - \xi) \\ & \text{subject to} && c_i(x) = 0, \quad i \in \mathcal{E} \\ & && c_i(x) \leq 0, \quad i \in \mathcal{I} \\ & && 0 \leq \xi \leq \mathbf{1}_n \quad \text{and} \quad \xi \circ x = 0. \end{aligned} \tag{10}$$

As an example of the problem (8), consider the misclassification minimization problem that arises from the literature in support-vector machines [7, 24]. For given scalars $(C, \gamma, \varepsilon) > 0$ and data

$(x^i, y_i) \in \mathbb{R}^{n+1}$ for $i = 1, \dots, m$, seek a pair $(w, b) \in \mathbb{R}^{n+1}$ to

$$\begin{aligned} \underset{(w,b)}{\text{minimize}} \quad & \underbrace{C \sum_{i=1}^m \max \{ |w^T x^i + b - y_i| - \varepsilon, 0 \}}_{\text{standard SVM objective}} + \frac{1}{2} w^T w + \\ & \underbrace{\gamma \sum_{i=1}^m \begin{cases} 1 & \text{if } |w^T x^i + b - y_i| > \varepsilon \\ 0 & \text{if } |w^T x^i + b - y_i| \leq \varepsilon. \end{cases}}_{\text{number of misclassified points}} \end{aligned}$$

Introducing the errors $e_i \triangleq \max \{ |w^T x^i + b - y_i| - \varepsilon, 0 \}$, which we assemble in the m -vector $e \triangleq (e_i)_{i=1}^m$, we can reformulate the above misclassification minimization problem in the form of (8):

$$\begin{aligned} \underset{(w,e,b)}{\text{minimize}} \quad & C \sum_{i=1}^m e_i + \frac{1}{2} w^T w + \gamma \|e\|_0 \\ \text{subject to} \quad & \begin{cases} e_i \geq w^T x^i + b - y_i - \varepsilon \\ e_i \geq -(w^T x^i + b - y_i) - \varepsilon \\ e_i \geq 0 \end{cases} \quad i = 1, \dots, m. \end{aligned}$$

The objective function contains the regression errors measured in two ways: size and count, weighted by the parameter C .

3 A Touch of Piecewise Theory

In practice, the problem (10) provides a computational platform for solving the problem (8). Thus it is important to understand the basic connections between these two problems. Due to the presence of the bilinear constraints: $\xi \circ x = 0$, (10) is a nonconvex program even if the original NLP (8) with $\gamma = 0$ is convex. The discussion in this section focuses on the half-complementarity formulation (10). To avoid repetition, results for the full-complementarity formulation of the problem (8) are stated without proof.

3.1 Constraint qualifications

As a first step towards understanding the connections between these problems, we begin with the discussion of some constraint qualifications. Let x be a feasible solution to (8). We wish to postulate a constraint qualification at x with respect to (8) under which the constant rank constraint qualification (CRCQ) [13, page 262] [20] will hold for the constraints of (10) at (x, ξ) , where ξ is defined by (4). We introduce several important index sets associated with x . Let

$$\mathcal{N}(x) \triangleq \{i \mid x_i \neq 0\} \quad \text{and} \quad \mathcal{A}(x) \triangleq \{i \in \mathcal{I} \mid c_i(x) = 0\}.$$

Also let $\mathcal{N}(x)^c$ be the complement of $\mathcal{N}(x)$ in $\{1, \dots, n\}$. The gradients of the active constraints in (10) at the pair (x, ξ) are of several kinds:

$$\left\{ \begin{pmatrix} \nabla c_i(x) \\ 0 \end{pmatrix} : i \in \mathcal{E} \cup \mathcal{A}(x) \right\}, \quad \left\{ - \begin{pmatrix} 0 \\ e_i \end{pmatrix} : i \in \mathcal{N}(x) \right\}, \quad \left\{ \begin{pmatrix} 0 \\ e_i \end{pmatrix} : i \in \mathcal{N}(x)^c \right\}$$

and $\underbrace{\left\{ x_i \begin{pmatrix} 0 \\ e_i \end{pmatrix} : i \in \mathcal{N}(x) \right\}, \left\{ \xi_i \begin{pmatrix} e_i \\ 0 \end{pmatrix} : i \in \mathcal{N}(x)^c \right\}}_{\text{gradients of the equality constraint } \xi \circ x = 0}$

where e_i is the n -vector of zeros except for a 1 in the i th position. We assume that for every index set $\alpha \subseteq \mathcal{A}(x)$ the family of vectors

$$\left\{ \left(\frac{\partial c_i(x')}{\partial x_j} \right)_{j \in \mathcal{N}(x)} : i \in \mathcal{E} \cup \alpha \right\} \quad (11)$$

has the same rank for all vectors x' sufficiently close to x that are also feasible to (8). Each vector (11) is a subvector of the gradient vector $\nabla c_i(x')$ with the partial derivatives $\partial c_i(x')/\partial x_j$ for $j \in \mathcal{N}(x)^c$ removed. If this assumption holds at x , the CRCQ is valid for the constraints of the problem (10) at the pair (x, ξ) . To show this, it suffices to verify that for any index sets $\alpha \subseteq \mathcal{A}(x)$, $\beta_1, \gamma_1 \subseteq \mathcal{N}(x)$ and $\beta_2, \gamma_2 \subseteq \mathcal{N}(x)^c$, the family of vectors:

$$\left\{ \begin{pmatrix} \nabla c_i(x') \\ 0 \end{pmatrix} : i \in \mathcal{E} \cup \alpha \right\}, \quad \left\{ - \begin{pmatrix} 0 \\ e_i \end{pmatrix} : i \in \beta_1 \right\}, \quad \left\{ \begin{pmatrix} 0 \\ e_i \end{pmatrix} : i \in \beta_2 \right\}$$

and $\underbrace{\left\{ x'_i \begin{pmatrix} 0 \\ e_i \end{pmatrix} : i \in \gamma_1 \right\}, \left\{ \xi'_i \begin{pmatrix} e_i \\ 0 \end{pmatrix} : i \in \gamma_2 \right\}}_{\substack{\text{gradients of the equality constraint } \xi \circ x = 0 \\ \text{evaluated at } (x', \xi')}} \quad (12)$

has the same rank for all pairs (x', ξ') sufficiently close to the given pair (x, ξ) that are also feasible to (10). Consider such a pair (x', ξ') . We must have $\mathcal{N}(x) \subseteq \mathcal{N}(x')$; moreover, if $i \in \mathcal{N}(x)^c$, then $\xi_i = 1$; hence $\xi'_i > 0$. By complementarity, it follows that $i \in \mathcal{N}(x')^c$. Consequently, $\mathcal{N}(x) = \mathcal{N}(x')$. This is sufficient to establish the rank invariance of the vectors (12) when the pair (x', ξ') varies near (x, ξ) .

3.2 Global optimality

For any index set $\mathcal{J} \subseteq \{1, \dots, n\}$ with complement \mathcal{J}^c , consider the nonlinear program

$$\begin{aligned} & \underset{x}{\text{minimize}} && f(x) \\ & \text{subject to} && c_i(x) = 0, \quad i \in \mathcal{E} \\ & && c_i(x) \leq 0, \quad i \in \mathcal{I} \\ & \text{and} && x_i = 0, \quad i \in \mathcal{J}, \end{aligned} \quad (13)$$

which may be thought of as a “piece” of (8) in the sense of piecewise programming. Indeed, provided that (8) is feasible, we have

$$\boxed{\infty > \text{minimum of (8)} = \min_{\mathcal{J}} \{ \text{minimum of (13)} + \gamma |\mathcal{J}^c| \} \geq -\infty,} \quad (14)$$

where the value of $-\infty$ is allowed in both the left- and right-hand sides. [We adopt the convention that the minimum value of an infeasible optimization problem is taken to be ∞ .] To prove (14), we note that the left-hand minimum is always an upper bound of the right-hand minimum. Indeed, for any feasible x of (8), we have, with $\mathcal{J} = \mathcal{N}(x)^c$,

$$f(x) + \gamma \|x\|_0 = f(x) + \gamma |\mathcal{J}^c| \geq \{ \text{minimum of (13) with } \mathcal{J} = \mathcal{N}(x)^c \} + \gamma |\mathcal{J}^c|.$$

This bound establishes the equality of the two minima in (14) when the left-hand minimum is equal to $-\infty$. A moment's thought shows that these two minima are also equal when the right-hand minimum is equal to $-\infty$. Thus it remains to consider the case where both the left and right minima are finite. Let $x^{\mathcal{J}}$ be an optimal solution of (13) that attains the right-hand minimum in (14). We have

$$\text{minimum of (8)} \leq f(x^{\mathcal{J}}) + \gamma \|x^{\mathcal{J}}\|_0 \leq f(x^{\mathcal{J}}) + \gamma |\mathcal{J}^c|.$$

This completes the proof of the claim of the equality (14). An immediate consequence of this equality is that it provides sufficient conditions under which the discontinuous ℓ_0 -minimization problem (8) has an optimal solution. For instance, if the objective function f and the constraint functions c_i are all continuous and that f is coercive on the feasible set of (8), then the ℓ_0 -problem has an optimal solution.

3.3 KKT conditions and local optimality

With the CRCQ in place, it follows that the Karush-Kuhn-Tucker (KKT) conditions are necessary for a pair (x, ξ) to be optimal to (10), provided that ξ is defined by (4). Letting λ , η , and μ be the multipliers to the constraints of (10), the KKT conditions of this problem are:

$$\begin{aligned} 0 &= \nabla f(x) + \sum_{i \in \mathcal{E} \cup \mathcal{I}} \lambda_i \nabla c_i(x) + \mu \circ \xi \\ 0 &\leq \xi \perp -\gamma + \mu \circ x + \eta \geq 0 \\ 0 &\leq \eta \perp \mathbf{1}_n - \xi \geq 0 \\ 0 &= c_i(x), \quad i \in \mathcal{E} \\ 0 &\leq \lambda_i \perp c_i(x) \leq 0, \quad i \in \mathcal{I} \\ 0 &= \xi \circ x. \end{aligned} \tag{15}$$

Letting λ denote the multipliers of the functional constraints in (13), we can write the KKT conditions of (13) as

$$\begin{aligned} 0 &= \frac{\partial f(x)}{\partial x_j} + \sum_{i \in \mathcal{E} \cup \mathcal{I}} \lambda_i \frac{\partial c_i(x)}{\partial x_j}, \quad j \in \mathcal{J}^c \\ 0 &= c_i(x), \quad i \in \mathcal{E} \\ 0 &\leq \lambda_i \perp c_i(x) \leq 0, \quad i \in \mathcal{I} \\ 0 &= x_j, \quad j \in \mathcal{J}. \end{aligned} \tag{16}$$

We have the following result connecting the KKT systems (15) and (16), which can be contrasted with the equality (14) that deals with the global minima of these two problems.

Proposition 1 Let x be a feasible solution of (8) and let ξ be defined by (4). The following three statements hold for any positive vector γ :

- (a) If (x, ξ) is a KKT point of (10) with multipliers λ , μ , and η , then x is a KKT point of (13) for any \mathcal{J} satisfying $\mathcal{N}(x) \subseteq \mathcal{J}^c \subseteq \mathcal{N}(x) \cup \{j \mid \mu_j = 0\}$.
- (b) Conversely, if x is a KKT point of (13) for some \mathcal{J} , then (x, ξ) is a KKT point of (10).
- (c) If x is a local minimum of (13) for some \mathcal{J} , then (x, ξ) is a local minimum of (10).

Proof. To prove (a), it suffices to note that for such an index set \mathcal{J} , we must have $\mathcal{J} \subseteq \mathcal{N}(x)^c$; moreover, if $j \in \mathcal{J}^c$, then either $\mu_j = 0$ or $\xi_j = 0$. To prove part (b), it suffices to define the multipliers μ and η . An index set \mathcal{J} for which (16) holds must be a subset of $\mathcal{N}(x)^c$ so that $\mathcal{N}(x) \subseteq \mathcal{J}^c$. Let

$$\mu_j \triangleq \begin{cases} \frac{\gamma_j}{x_j} & \text{if } j \in \mathcal{N}(x) \\ -\left[\frac{\partial f(x)}{\partial x_j} + \sum_{i \in \mathcal{E} \cup \mathcal{I}} \lambda_i \frac{\partial c_i(x)}{\partial x_j} \right] & \text{if } j \in \mathcal{J} \\ 0 & \text{if } j \in \mathcal{J}^c \cap \mathcal{N}(x)^c \end{cases} \quad \eta_j \triangleq \begin{cases} 0 & \text{if } j \in \mathcal{N}(x) \\ \gamma_j & \text{if } j \in \mathcal{J} \\ \gamma_j & \text{if } j \in \mathcal{J}^c \cap \mathcal{N}(x)^c. \end{cases}$$

It is not difficult to verify that the KKT conditions (15) hold at the triple $(x, \xi, \lambda, \mu, \eta)$.

Finally, to prove (c), let (x', ξ') be a feasible pair of (10) sufficiently close to (x, ξ) . Since $x_j = 0$ for all $j \in \mathcal{J}$, it follows that $\xi_j = 1$ for all such j . Since ξ'_j is sufficiently close to ξ_j , we deduce $\xi'_j > 0$; hence $x'_j = 0$ by complementarity. Thus x' is feasible to (13). Moreover, if $x_j \neq 0$, then $x'_j \neq 0$; hence $\xi'_j = 0$.

$$\begin{aligned} & f(x') + \gamma^T (\mathbf{1}_n - \xi') \\ & \geq f(x) + \sum_{j: x_j=0} \gamma_j (1 - \xi'_j) + \sum_{j: x_j \neq 0} \gamma_j (1 - \xi'_j) \\ & \geq f(x) + \sum_{j: x_j=0} \gamma_j (1 - \xi_j) + \sum_{j: x_j \neq 0} \gamma_j (1 - \xi_j) = f(x) + \gamma^T (\mathbf{1}_n - \xi), \end{aligned}$$

establishing the desired conclusion. \square

The next proposition states a similar result for the inner-product and componentwise reformulations of the full-complementarity formulation (9).

Proposition 2 Let x be a feasible solution of (8), let ξ be defined by (4), and $x^\pm \triangleq \max\{0, \pm x\}$. Then the following three statements hold for any positive vector γ :

- (a) If (x, x^\pm, ξ) is a KKT point of

$$\begin{aligned} & \underset{x, x^\pm, \xi}{\text{minimize}} && f(x) + \gamma^T (\mathbf{1}_n - \xi) \\ & \text{subject to} && c_i(x) = 0, \quad i \in \mathcal{E} \\ & && c_i(x) \leq 0, \quad i \in \mathcal{I} \\ & && x = x^+ - x^-, \quad \xi^T (x^+ + x^-) \leq 0, \quad (x^+)^T (x^-) \leq 0 \\ & && 0 \leq \xi \leq \mathbf{1}_n \quad \text{and} \quad x^\pm \geq 0, \end{aligned} \tag{17}$$

or of

$$\begin{aligned}
& \underset{x, x^\pm, \xi}{\text{minimize}} && f(x) + \gamma^T (\mathbf{1}_n - \xi) \\
& \text{subject to} && c_i(x) = 0, \quad i \in \mathcal{E} \\
& && c_i(x) \leq 0, \quad i \in \mathcal{I} \\
& && x = x^+ - x^-, \quad \xi \circ (x^+ + x^-) \leq 0, \quad (x^+) \circ (x^-) \leq 0 \\
& && 0 \leq \xi \leq \mathbf{1}_n \quad \text{and} \quad x^\pm \geq 0
\end{aligned} \tag{18}$$

then x is a KKT point of (13) for any \mathcal{J} satisfying $\mathcal{N}(x) \subseteq \mathcal{J}^c \subseteq \mathcal{N}(x) \cup \{j \mid \mu_j = 0\}$.

(b) Conversely, if x is a KKT point of (13) for some \mathcal{J} , then (x, x^\pm, ξ) is a KKT point of (17) and (18).

(c) If x is a local minimum of (13) for some \mathcal{J} , then (x, x^\pm, ξ) is a local minimum of (17) and (18). \square

We now look at the second-order optimality conditions. In the proposition below, we examine the sufficient conditions; an analogous result can be derived for the second-order necessary conditions in a similar manner.

Proposition 3 Let (x, ξ) be a point so that (x, ξ) is a KKT point of (10) with multipliers λ , μ , and η , and so that x is a KKT point of (13) for any \mathcal{J} satisfying $\mathcal{N}(x) \subseteq \mathcal{J}^c \subseteq \mathcal{N}(x) \cup \{j \mid \mu_j = 0\}$. If the second-order sufficient condition holds at (x, ξ) of (13), then it holds at x of (10). In addition, if $\mathcal{J}^c = \mathcal{N}(x)$, the converse holds.

Proof. The second-order conditions examine directions d in the critical cone, i.e., those directions that satisfy the linearization of each equality constraint and active inequality constraints, and that keep active the linearization of any inequality constraint with a positive multiplier. From the KKT conditions (15) of (10), if $x_j = 0$ for some $j \in \{1, \dots, n\}$ then $\xi_j = 1$ and the corresponding multiplier $\eta_j > 0$. If $x_j \neq 0$ then the linearization of the complementarity constraint restricts the direction. Thus, all directions $d = (d_x, d_\xi)$ in the critical cone satisfy

$$d_\xi = 0.$$

Therefore, we only need to consider the x part of the Hessian,

$$H = \nabla^2 f(x) + \sum_{i \in \mathcal{E} \cup \mathcal{I}} \lambda_i \nabla^2 c_i(x),$$

for both (10) and (13). Let $D_1 \subseteq \mathbb{R}^n$ be the set of directions d_x that satisfy

$$\begin{aligned}
\nabla c_i(x)^T d_x &= 0, & i \in \mathcal{E} \\
\nabla c_i(x)^T d_x &= 0, & i \in \mathcal{I} \text{ with } \lambda_i > 0 \\
\nabla c_i(x)^T d_x &\leq 0, & i \in \mathcal{I} \text{ with } \lambda_i = 0
\end{aligned} \tag{19}$$

together with $d_x \circ \xi = 0$. The second-order sufficient condition for (10) holds at x if and only if $d^T H d > 0$ for all $d \in D_1$ with $d \neq 0$. Similarly, let $D_2 \subseteq \mathbb{R}^n$ be the set of directions that satisfy (19) together with $(d_x)_j = 0$ for all $j \in \mathcal{J}$. Then the second-order sufficient condition for (13) holds at x if and only if $d^T H d > 0$ for all $d \in D_2$ with $d \neq 0$.

To prove the first part of the claim, we need to show that $D_1 \subseteq D_2$. Let $d^1 \in D_1$. Since $\mathcal{J} \subseteq \mathcal{N}(x)^c$, we have $\xi_j = 1$ for all $j \in \mathcal{J}$. Because the direction satisfies the linearization of the half-complementarity constraint, it follows that $d_j^1 = 0$, which implies $d^1 \in D_2$. To prove the second part, let $d^2 \in D_2$. Since $\mathcal{N}(x)^c = \mathcal{J}$, we have $x_j = 0$ and hence $\xi_j = 1$ for all $j \in \mathcal{J}$. Further, $x_j \neq 0$ and $\xi_j = 0$ for all $j \in \mathcal{J}^c$. Thus $d^2 \in D_1$. \square

Again, similar relationships between the second-order optimality conditions for (13), (17), and (18) can be established, but we omit the formal statement for the sake of brevity.

Let us consider what these results imply for the simple model problem (1). It is easy to see that *any feasible point* x with $Ax \geq b$ and $Cx = d$ is a KKT point for (13) with $\mathcal{J} = \{j : x_j = 0\}$ and $f(x) = 0$. Propositions 1(b) and 3 then imply that x corresponds to a KKT point of (10) that satisfies the second-order necessary optimality conditions. In other words, finding a second-order necessary KKT point for (10) merely implies that we found a *feasible* point for (1), but this says nothing about its ℓ_0 -norm. We summarize this observation in the following corollary. An analogous result is true for the formulations (17) and (18).

Corollary 1 A vector \hat{x} is feasible to the system $Ax \geq b$ and $Cx = d$ if and only if (\hat{x}, ξ) , where $\mathbf{1}_n - \xi$ is the support vector of \hat{x} , is a KKT pair of the nonlinear program (10) with $\mathcal{J} = \{j : \hat{x}_j = 0\}$ that satisfies the second-order necessary optimality conditions. \square

A principal goal in this study is assessing the adequacy of (local) NLP solvers for solving ℓ_0 -minimization problems, such as (1) and (8), using the equivalent full- or half-complementarity reformulations. The results in this section cast a negative shadow on this approach. NLP solvers typically aim to find KKT points, ideally those that satisfy second-order optimality conditions. Propositions 1–3 establish that the reformulations for (8) may have an exponential number of points (one for each set $\mathcal{J} \subseteq \{1, \dots, n\}$ in (13)), to which the NLP solvers might be attracted to. In the particular case of the model problem (1), Corollary 1 paints an even more gloomy picture because *any* feasible point for (1) has the characteristics that an NLP solver looks for, and most of those points have sub-optimal objective function values. Interestingly, these theoretical reservations do not seem to materialize in practice. Our computational results attest that usually points of rather good objective values are returned by the NLP solvers. The discussions related to the relaxed formulations in Section 4 shed some light on this observation.

3.4 KKT conditions for the squared reformulation

To conclude the discussion of exact smooth reformulations of the ℓ_0 -problem, we examine the new squared complementarity reformulation. Transforming the complementarity conditions in (9) using

this approach, we obtain the NLP

$$\begin{aligned}
& \underset{x, x^\pm, \xi}{\text{minimize}} && f(x) + \gamma^T (\mathbf{1}_n - \xi) \\
& \text{subject to} && c_{\mathcal{E}}(x) = 0 && (\lambda_{\mathcal{E}}) \\
& && c_{\mathcal{I}}(x) \leq 0 && (\lambda_{\mathcal{I}}) \\
& && x - x^+ + x^- = 0 && (\lambda_x) \\
& && \phi(x^+, x^-, \xi) \triangleq [\xi^T(x^+ + x^-) + (x^+)^T(x^-)]^2 \leq 0 && (\mu) \\
& && \xi \leq \mathbf{1}_n && (\eta) \\
& \text{and} && x^\pm, \xi \geq 0.
\end{aligned} \tag{20}$$

Independently of a constraint qualification, we may write down the KKT conditions for this problem. These include

$$\begin{aligned}
0 \leq \xi & \perp -\gamma + [J_\xi \phi(x^+, x^-, \xi)]^T \mu + \eta \geq 0 \\
0 \leq \eta & \perp \mathbf{1}_n - \xi \geq 0,
\end{aligned} \tag{21}$$

where J_ξ denotes the (partial) Jacobian matrix with respect to the variable ξ . It is easy to see that for any feasible point, $\phi(x^+, x^-, \xi)$ as well as its first derivatives are zero. Therefore, (21) implies that $\xi = \mathbf{1}_n$ and $\eta = \gamma$. Furthermore, $\phi(x^+, x^-, \xi) = 0$ gives that $x^+ = x^- = 0$. Thus, we have proved

Proposition 4 The only KKT point of (20), if it exists, is $(x, x^+, x^-, \xi) = (0, 0, 0, \mathbf{1}_n)$. Hence, if $x = 0$ is not feasible to (20), then no KKT point exists for this problem. \square

Again, from a theoretical perspective, this result suggests that it is not a good idea to use an NLP algorithm to solve the ℓ_0 -problems (1) or (8) transformed by the squared reformulation. Nevertheless, our numerical experiments reported in Section 5 suggests otherwise. Indeed, the encouraging computational results are the primary reason for us to introduce this squared formulation in the first place.

4 Relaxed Formulations

As mentioned at the end of the introduction, in general, the exact reformulations of an MPCC result in NLPs that are not well-posed in the sense that the MFCQ does not hold at any feasible point. To overcome this shortcoming, relaxation schemes for MPCCs have been proposed [26, 10, 21], where the complementarity constraints are relaxed. The resulting NLPs have better properties, and the solution of the original MPCC is obtained by solving a sequence of relaxed problems, for which the relaxation parameter is driven to zero. In this section, we investigate the stationarity properties of relaxed reformulations for the ℓ_0 problem.

We introduce the following relaxation of the new half-complementarity formulation (10), which we

denote by $\text{NLP}(\varepsilon)$, for a given relaxation scalar $\varepsilon > 0$:

$$\begin{aligned}
& \underset{x, \xi}{\text{minimize}} && f(x) + \gamma^T (\mathbf{1}_n - \xi) \\
& \text{subject to} && c_i(x) = 0, \quad i \in \mathcal{E} && (\lambda_{\mathcal{E}}) \\
& && c_i(x) \leq 0, \quad i \in \mathcal{I} && (\lambda_{\mathcal{I}}) \\
& && \xi \leq \mathbf{1}_n && (\eta) \\
& && \xi \circ x \leq \varepsilon \mathbf{1}_n && (\mu^+) \\
& && -\xi \circ x \leq \varepsilon \mathbf{1}_n && (\mu^-) \\
& \text{and} && \xi \geq 0,
\end{aligned} \tag{22}$$

where λ , η , and μ^\pm are the associated multipliers for the respective constraints. The problem $\text{NLP}(0)$ is the original half-complementary formulation (10). Observations analogous to those in the following sections are valid for relaxations of the full complementarity formulations (17) and (18).

4.1 Convergence of KKT points for $\text{NLP}(\varepsilon)$

In this section we examine the limiting behavior of KKT points $x(\varepsilon)$ for $\text{NLP}(\varepsilon)$ as ε converges to zero. This is of interest because algorithms based on the sequential solution of the relaxation $\text{NLP}(\varepsilon)$ aim to compute limit points of $x(\varepsilon)$ [10, 21, 26]. However, our analysis also gives insight into the behavior of standard NLP solvers that are applied directly to one of the unrelaxed NLP-reformulations of the MPCC, such as (10), (17), and (18). For instance, some implementations of NLP solvers, such as the IPOPT solver [28] used for the numerical experiments in Section 5, relax all inequality and bound constraints by a small amount that is related to the convergence tolerance before solving the problem at hand. This modification is done in order to make the problem somewhat “nicer”; for example, a feasible problem is then guaranteed to have a nonempty relative interior of the feasible region. However, because this alteration is on the order of the user-specified convergence tolerance, it usually does not lead to solutions that are far away from solutions of the original unperturbed problem. In the current context this means that such an NLP code solves the relaxation $\text{NLP}(\varepsilon)$ even if the relaxation is not explicitly introduced by the user.

Before examining the limiting behavior of $x(\varepsilon)$ on a specific example, we consider the KKT conditions for $\text{NLP}(\varepsilon)$:

$$0 = \nabla f(x) + \sum_{i \in \mathcal{E} \cup \mathcal{I}} \lambda_i \nabla c_i(x) + (\mu^+ - \mu^-) \circ \xi \tag{23a}$$

$$0 \leq \xi \perp -\gamma + (\mu^+ - \mu^-) \circ x + \eta \geq 0 \tag{23b}$$

$$0 \leq \eta \perp \mathbf{1}_n - \xi \geq 0 \tag{23c}$$

$$0 = c_i(x), \quad i \in \mathcal{E} \tag{23d}$$

$$0 \leq \lambda_i \perp c_i(x) \leq 0, \quad i \in \mathcal{I} \tag{23e}$$

$$0 \leq \mu^+ \perp \varepsilon \mathbf{1}_n - \xi \circ x \geq 0 \tag{23f}$$

$$0 \leq \mu^- \perp \varepsilon \mathbf{1}_n + \xi \circ x \geq 0. \tag{23g}$$

For a KKT point (x, ξ) , these conditions imply the following:

• For $j \in \{1, \dots, n\}$ with $x_j > \varepsilon > 0$, by (23g) we have $\mu_j^- = 0$, and by (23f) we have $\xi_j < 1$. It follows from (23c) that $\eta_j = 0$ and then (23b) and (23a) give the relationships:

$$\xi_j = \frac{\varepsilon}{x_j} < 1, \quad \eta_j = 0, \quad \mu_j^+ = \frac{\gamma_j}{x_j}, \quad \mu_j^- = 0, \quad \frac{\partial f(x)}{\partial x_j} + \sum_{i \in \mathcal{E} \cup \mathcal{I}} \lambda_i \frac{\partial c_i(x)}{\partial x_j} = -\frac{\varepsilon \gamma_j}{x_j^2} < 0. \quad (24)$$

• For $j \in \{1, \dots, n\}$ with $x_j < -\varepsilon < 0$, by (23f) we have $\mu_j^+ = 0$, and by (23g) we have $\xi_j < 1$. It follows from (23c) that $\eta_j = 0$ and then (23b) and (23a) give the relationships:

$$\xi_j = \frac{\varepsilon}{-x_j} < 1, \quad \eta_j = 0, \quad \mu_j^+ = 0, \quad \mu_j^- = \frac{\gamma_j}{-x_j}, \quad \frac{\partial f(x)}{\partial x_j} + \sum_{i \in \mathcal{E} \cup \mathcal{I}} \lambda_i \frac{\partial c_i(x)}{\partial x_j} = \frac{\varepsilon \gamma_j}{x_j^2} > 0. \quad (25)$$

• For $j \in \{1, \dots, n\}$ with $-\varepsilon < x_j < \varepsilon$, (23f) and (23g) give $\mu_j^+ = \mu_j^- = 0$. Then (23b) implies $\eta_j > 0$, giving $\xi_j = 1$ by (23c) and so $\eta_j = \gamma_j$ by (23b). Together with (23a), this overall implies

$$\xi_j = 1, \quad \eta_j = \gamma_j, \quad \mu_j^+ = 0, \quad \mu_j^- = 0, \quad \frac{\partial f(x)}{\partial x_j} + \sum_{i \in \mathcal{E} \cup \mathcal{I}} \lambda_i \frac{\partial c_i(x)}{\partial x_j} = 0. \quad (26)$$

• For $j \in \{1, \dots, n\}$ with $x_j = \varepsilon$ we have $\mu_j^- = 0$ from (23g). Also, we must have $\xi_j = 1$ since otherwise $\eta_j = 0$ from (23c) and $\mu_j^+ = 0$ from (23f), which would then violate (23b). Thus from (23b) and (23a) we have

$$\xi_j = 1, \quad \eta_j = \gamma_j - \varepsilon \mu_j^+, \quad 0 \leq \mu_j^+ \leq \frac{\gamma_j}{\varepsilon}, \quad \mu_j^- = 0, \quad \frac{\partial f(x)}{\partial x_j} + \sum_{i \in \mathcal{E} \cup \mathcal{I}} \lambda_i \frac{\partial c_i(x)}{\partial x_j} + \mu_j^+ = 0. \quad (27)$$

• For $j \in \{1, \dots, n\}$ with $x_j = -\varepsilon$ we have $\mu_j^+ = 0$ from (23f). Also, we must have $\xi_j = 1$ since otherwise $\eta_j = 0$ from (23c) and $\mu_j^- = 0$ from (23g), which would then violate (23b). Thus from (23b) and (23a) we have

$$\xi_j = 1, \quad \eta_j = \gamma_j - \varepsilon \mu_j^-, \quad \mu_j^+ = 0, \quad 0 \leq \mu_j^- \leq \frac{\gamma_j}{\varepsilon}, \quad \frac{\partial f(x)}{\partial x_j} + \sum_{i \in \mathcal{E} \cup \mathcal{I}} \lambda_i \frac{\partial c_i(x)}{\partial x_j} - \mu_j^- = 0. \quad (28)$$

We now use these observations to characterize the limit points of KKT points $x(\varepsilon)$ for NLP(ε) in a particular case. We consider an instance of the model problem (1) with the simple set of linear constraints

$$x_1 + x_2 + x_3 \geq 1$$

$$x_1, x_2, x_3 \geq 0.$$

The relaxation of the corresponding half-complementarity formulation is

$$\begin{aligned} & \underset{x, \xi}{\text{minimize}} && \mathbf{1}_3^T (\mathbf{1}_3 - \xi) \\ & \text{subject to} && 1 - x_1 - x_2 - x_3 \leq 0, && \lambda_0 \\ & && -x_i \leq 0, && \lambda_i \quad \text{for } i = 1, 2, 3 \\ & && \xi \leq \mathbf{1}_3 && \eta \\ & && \xi \circ x \leq \varepsilon \mathbf{1}_3 && \mu^+ \\ & \text{and} && \xi \geq 0, \end{aligned} \quad (29)$$

which is a special case of (22) where $f(x) = 0$, $\gamma = \mathbf{1}_3$, $c_0(x) = 1 - x_1 - x_2 - x_3$, and $c_i(x) = -x_i$ for $i = 1, 2, 3$. Because here $\xi, x \geq 0$ for any feasible point, the constraint corresponding to μ^- in (22) can never be active, and we may ignore this constraint without loss of generality.

The following proposition shows that there are exactly seven limit points of $x(\varepsilon)$ as ε converges to zero, and that the KKT points converging to non-global solutions of the unrelaxed half-complementarity formulation (10) are local maximizers of (29).

Proposition 5 The following statements are true.

(a) The set of limit points of KKT points for (29) as $\varepsilon \rightarrow 0$ is exactly

$$\left\{ \left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3} \right), \left(0, \frac{1}{2}, \frac{1}{2} \right), \left(\frac{1}{2}, 0, \frac{1}{2} \right), \left(\frac{1}{2}, \frac{1}{2}, 0 \right), (0, 0, 1), (0, 1, 0), (1, 0, 0) \right\}.$$

(b) The KKT points $x(\varepsilon)$ of (29) that converge to

$$\hat{x} \in \left\{ \left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3} \right), \left(0, \frac{1}{2}, \frac{1}{2} \right), \left(\frac{1}{2}, 0, \frac{1}{2} \right), \left(\frac{1}{2}, \frac{1}{2}, 0 \right) \right\}$$

as $\varepsilon \rightarrow 0$ are strict local maximizers.

Proof. Part (a) can be proven by specializing the conclusions (24)–(28) to (29). Let $x(\varepsilon)$ be a sequence of KKT points for (29) that converges to some limit point \hat{x} as $\varepsilon \rightarrow 0$, and $\xi(\varepsilon)$, $\eta(\varepsilon)$, $\mu^+(\varepsilon)$, $\lambda(\varepsilon)$ the corresponding quantities in the KKT conditions (23).

Due to the first constraint in (29), $x(\varepsilon)$ has at least one component $x_j(\varepsilon) \geq \frac{1}{3} > \varepsilon$ if ε is sufficiently small. For such j we know that $\lambda_j(\varepsilon) = 0$ due to $x_j(\varepsilon) > 0$. Moreover, based on (24) we have

$$\frac{\partial f(x)}{\partial x_j} + \sum_{i \in \mathcal{E} \cup \mathcal{I}} \lambda_i \frac{\partial c_i(x)}{\partial x_j} = -\lambda_0(\varepsilon) - \lambda_i(\varepsilon) = -\frac{\varepsilon}{x_j(\varepsilon)^2}$$

so that

$$\lambda_0(\varepsilon) = \frac{\varepsilon}{x_j(\varepsilon)^2} > 0 \tag{30}$$

for such j . Hence the constraint corresponding to $\lambda_0(\varepsilon)$ is active and we have

$$x_1(\varepsilon) + x_2(\varepsilon) + x_3(\varepsilon) = 1. \tag{31}$$

We consider three cases, depending on the nonzero structure of \hat{x} .

Case 1: Suppose $\hat{x}_j > 0$ for all $j = 1, 2, 3$. For sufficiently small ε we have $x_j(\varepsilon) > \varepsilon$ and therefore $\lambda_j(\varepsilon) = 0$ for all $j = 1, 2, 3$. Then, based on (24), we know that (30) holds for all $j = 1, 2, 3$, so that $x_1(\varepsilon) = x_2(\varepsilon) = x_3(\varepsilon)$. Using (31) and (24) we see that the KKT quantities must satisfy

$$\begin{aligned} x(\varepsilon) &= \left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3} \right), & \xi(\varepsilon) &= (3\varepsilon, 3\varepsilon, 3\varepsilon), & \lambda(\varepsilon) &= (9\varepsilon, 0, 0, 0) \\ \mu^+(\varepsilon) &= (3, 3, 3), & \eta(\varepsilon) &= (0, 0, 0). \end{aligned} \tag{32}$$

This shows that $\lim_{\varepsilon \rightarrow 0} x(\varepsilon) = \left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3} \right)$ is the only potential limit point with the nonzero structure considered in this case. Conversely, it is easy to verify that the quantities in (32) satisfy the KKT conditions (23), so that $\hat{x} = \left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3} \right)$ is indeed a limit point.

Case 2: Suppose $\hat{x} = (0, \hat{x}_2, \hat{x}_3)$ with $\hat{x}_2, \hat{x}_3 > 0$ (similarly for points with the structure $(\hat{x}_1, 0, \hat{x}_3)$ and $(\hat{x}_1, \hat{x}_2, 0)$). First, we note that $x_1(\varepsilon) \leq \varepsilon$ for small ε , since otherwise as in Case 1 we have $x_1(\varepsilon) = \frac{1}{3} \not\rightarrow 0$. If $x_1(\varepsilon) < \varepsilon$, then (26) would dictate that

$$\frac{\partial f(x)}{\partial x_j} + \sum_{i \in \mathcal{E} \cup \mathcal{I}} \lambda_i \frac{\partial c_i(x)}{\partial x_j} = -\lambda_0(\varepsilon) - \lambda_1(\varepsilon) = 0.$$

which implies that $\lambda_0(\varepsilon) \leq 0$ because $\lambda_1(\varepsilon) \geq 0$. However, this contradicts with (30). Therefore, we must have $x_1(\varepsilon) = \varepsilon$. Because $x_2(\varepsilon) > \varepsilon$ and $x_3(\varepsilon) > \varepsilon$ for small ε , (30) holds for $j = 2, 3$, and with (31) we obtain $x_2(\varepsilon) = x_3(\varepsilon) = (1 - \varepsilon)/2$. Based on (24) and (27), the KKT quantities in this case have to satisfy

$$\begin{aligned} x(\varepsilon) &= \left(\varepsilon, \frac{1 - \varepsilon}{2}, \frac{1 - \varepsilon}{2} \right), & \xi(\varepsilon) &= \left(1, \frac{2\varepsilon}{1 - \varepsilon}, \frac{2\varepsilon}{1 - \varepsilon} \right), \\ \lambda(\varepsilon) &= \left(\frac{4\varepsilon}{(1 - \varepsilon)^2}, 0, 0, 0 \right), & \mu^+(\varepsilon) &= \left(\frac{4\varepsilon}{(1 - \varepsilon)^2}, \frac{2}{(1 - \varepsilon)}, \frac{2}{(1 - \varepsilon)} \right), \\ \eta(\varepsilon) &= \left(1 - \frac{4\varepsilon^2}{(1 - \varepsilon)^2}, 0, 0 \right). \end{aligned} \quad (33)$$

This shows that $\lim_{\varepsilon \rightarrow 0} x(\varepsilon) = (0, \frac{1}{2}, \frac{1}{2})$ is the only potential limit point with the nonzero structure considered in this case. Conversely, it is easy to verify that the quantities in (33) satisfy the KKT conditions (23), so that $\hat{x} = (0, \frac{1}{2}, \frac{1}{2})$ is indeed a limit point.

Case 3: Suppose $\hat{x} = (0, 0, \hat{x}_3)$ with $\hat{x}_3 > 0$ (similarly for $\hat{x} = (\hat{x}_1, 0, 0)$ and $\hat{x} = (0, \hat{x}_2, 0)$). Similar arguments as those in Case 2 lead to the following KKT quantities:

$$\begin{aligned} x(\varepsilon) &= (\varepsilon, \varepsilon, 1 - 2\varepsilon), & \xi(\varepsilon) &= \left(1, 1, \frac{\varepsilon}{1 - 2\varepsilon} \right), \\ \lambda(\varepsilon) &= \left(\frac{\varepsilon}{(1 - 2\varepsilon)^2}, 0, 0, 0 \right), & \mu^+(\varepsilon) &= \left(\frac{\varepsilon}{(1 - 2\varepsilon)^2}, \frac{\varepsilon}{(1 - 2\varepsilon)^2}, \frac{1}{(1 - 2\varepsilon)} \right), \\ \eta(\varepsilon) &= \left(\frac{1 + 4\varepsilon - 3\varepsilon^2}{(1 - 2\varepsilon)^2}, \frac{1 + 4\varepsilon - 3\varepsilon^2}{(1 - 2\varepsilon)^2}, 0 \right) \end{aligned} \quad (34)$$

with $x(\varepsilon) \rightarrow (0, 0, 1)$.

To prove part (b), let $x(\varepsilon)$ be a KKT point for (22) for small ε . From part (a) we know that this corresponds to one of the cases (32)–(34).

The only nonlinear function in (22) is the constraint involving “ $x \circ \xi = 0$ ”, so that the Hessian of the Lagrangian for (22) at the KKT point $(x(\varepsilon), \xi(\varepsilon))$ is

$$H = \begin{bmatrix} 0 & D(\varepsilon) \\ D(\varepsilon) & 0 \end{bmatrix}$$

with $D(\varepsilon) = \text{diag}(\mu^+(\varepsilon))$. In all three cases (32)–(34) above we have $\mu^+(\varepsilon) > 0$, so that all nonlinear constraints are active. Assuming the variable order (x, ξ) , the Jacobian of those constraints at the KKT point is $J_{\text{nonlin}} = [\Xi(\varepsilon) \quad X(\varepsilon)]$ with $X = \text{diag}(x(\varepsilon))$ and $\Xi = \text{diag}(\xi(\varepsilon))$, the columns of

$$Z = \begin{bmatrix} X(\varepsilon) \\ -\Xi(\varepsilon) \end{bmatrix}$$

are a basis of the null space $\text{Null}(J_{\text{nonlin}})$ of J_{nonlin} , and the reduced Hessian of the Lagrangian with respect to J_{nonlin} is

$$Z^T H Z = -2D(\varepsilon)X(\varepsilon)\Xi(\varepsilon) \quad (35)$$

and negative definite.

In both cases (32)–(33), the inequality constraint corresponding to λ_0 is active, and in (33) the constraint “ $\xi_1 \leq 1$ ” is active. Together with the nonlinear constraint, the total number of active constraints is 4 and 5, respectively, and the null space of the Jacobian J_{act} of the active constraints has a dimension of at least 1. Since $\text{Null}(J_{\text{act}}) \subseteq \text{Null}(J_{\text{nonlin}})$, we have with (35) that the reduced Hessian with respect to J_{act} is negative definite. Because strict complementarity holds, this shows that the KKT points in (32)–(33) are strict local maximizers. \square

For this particular example, this result is as strong as we could hope for. There are only a few KKT points for each sufficiently small ε with an equal number of limit points as $\varepsilon \downarrow 0$. This is in stark contrast to Corollary 1, which does not differentiate between *any feasible* point. In addition, if the relaxation $\text{NLP}(\varepsilon)$ is solved by an NLP solver that is able to escape local maximizers, the point returned by the solver for small ε is close to a global minimizer of the original ℓ_0 -problem (1). While this discussion considered only the particular instance (29), it may point the way to further analysis that might explain the encouraging results observed in our numerical experiments.

4.2 Convergence of global minimizers for $\text{NLP}(\varepsilon)$

Finally, we examine the convergence of sequences of global minimizers to $\text{NLP}(\varepsilon)$ as $\varepsilon \downarrow 0$. We first give a sufficient condition under which such minimizers are bounded.

Proposition 6 Assume that the function $f(x)$ is coercive. Then there exists a positive scalar σ such that for any $\varepsilon > 0$ we have $\|x\| \leq \sigma$ for any globally optimal solution (x, ξ) of $\text{NLP}(\varepsilon)$.

Proof. Without loss of generality we may assume that the constraints of (8) are consistent. Thus, by the coercivity of f , the NLP

$$\begin{aligned} & \underset{x}{\text{minimize}} && f(x) \\ & \text{subject to} && c_i(x) = 0, \quad i \in \mathcal{E} \\ & \text{and} && c_i(x) \leq 0, \quad i \in \mathcal{I}, \end{aligned}$$

has a global minimizer, which we denote x^* . Let ξ^* be given by (4) with $x = x^*$. Then (x^*, ξ^*) is a feasible point of $\text{NLP}(\varepsilon)$ for any $\varepsilon > 0$. Since $f(x)$ is coercive, there exists $\sigma > 0$ such that $f(x) > f(x^*) + \gamma^T(\mathbf{1}_n - \xi^*)$ if $\|x\| > \sigma$. Thus, for any $\xi \in [0, 1]^n$, we have

$$f(x) + \gamma^T(\mathbf{1}_n - \xi) \geq f(x) > f(x^*) + \gamma^T(\mathbf{1}_n - \xi^*);$$

therefore, (x, ξ) is not a globally optimal solution to (22). \square

The next result shows that any limit point of a sequence of global minimizers to (22) as $\varepsilon \downarrow 0$ must be a global minimizer of (8). While this result is mainly of theoretical interest (since in practice, it is not possible to solve the problems (22) to global optimality), it nevertheless sheds light on the convergence of a common approach for approximating the complementarity constraints as smooth nonlinear programs.

Proposition 7 Let $\{\varepsilon_k\}$ be a sequence of positive scalars converging to 0. For each k , let (x^k, ξ^k) be a global minimizer of $\text{NLP}(\varepsilon_k)$. Then, any limit point of the sequence $\{x^k\}$ is a globally optimal solution of (8).

Proof. Without loss of generality, we may let x^* be the limit of the sequence $\{x^k\}$. Clearly, x^* is feasible to (8). We derive a contradiction by assuming that x^* is not a globally optimal solution of (8). Let $g(x) \triangleq f(x) + \gamma^T(\mathbf{1}_n - \xi(x))$ with $\xi(x)$ given by (4). Suppose \hat{x} is a point in the feasible region of (8) with $g(x^*) - g(\hat{x}) \triangleq \alpha > 0$ for some α .

We first consider the case that $x^* = 0$, so $g(x^*) = f(x^*)$. From the convergence of $\{x^k\}$, there exists k_0 so that $f(x^k) \geq f(x^*) - 0.5\alpha$ for all $k > k_0$. The value of \hat{x} in (8) is no larger than $g(\hat{x})$, which is better than the value of x^k by at least 0.5α , contradicting the assumption that x^k is optimal for (8).

Now we consider the case $x^* \neq 0$. Define $x_{\min}^* \triangleq \frac{1}{2} \min\{|x_j^*| : |x_j^*| > 0, j = 1, \dots, n\}$. Since the sequence $\{x^k\}$ is convergent to x^* , a positive integer k_1 exists such that $\|x^k - x^*\|_\infty < x_{\min}^*$ for all $k > k_1$. For all $k > k_1$, it is easy to verify that we have $|x_j^k| \geq x_{\min}^*$ if $|x_j^*| > 0$. As $\varepsilon_k \downarrow 0$, a positive integer $k_2 \geq k_1$ exists such that $\varepsilon_{k_2} < \frac{\alpha x_{\min}^*}{\gamma^T \mathbf{1}_n}$. Then the global optimal solution (x^{k_2}, ξ^{k_2}) to $\text{NLP}(\varepsilon_{k_2})$ satisfies:

$$\begin{aligned} f(x^{k_2}) + \gamma^T(\mathbf{1}_n - \xi^{k_2}) &\geq f(x^{k_2}) + \sum_{j: |x_j^{k_2}| > 0} \gamma_j \left(1 - \frac{\varepsilon_{k_2}}{|x_j^{k_2}|}\right) \\ &\geq g(x^{k_2}) - \sum_{j: |x_j^{k_2}| > 0} \frac{\gamma_j \varepsilon_{k_2}}{x_{\min}^*} \geq g(x^*) - \frac{\varepsilon_{k_2} \gamma^T \mathbf{1}_n}{x_{\min}^*} > g(x^*) - \alpha = g(\hat{x}). \end{aligned}$$

The point $(\hat{x}, \xi(\hat{x}))$ is feasible in $\text{NLP}(\varepsilon_{k_2})$, so we have a contradiction to the fact that (x^{k_2}, ξ^{k_2}) is the global minimum to $\text{NLP}(\varepsilon_{k_2})$. \square

To close this section, we derive a specialized result for the linear feasibility problem (1). Sharpening Proposition 6, this result shows that the optimal solutions of the following ε -relaxed NLP:

$$\begin{aligned} &\underset{x, \xi}{\text{minimize}} && \mathbf{1}_n^T(\mathbf{1}_n - \xi) \\ &\text{subject to} && Ax \geq b; \quad Cx = d \\ & && \xi \circ x \leq \varepsilon \mathbf{1}_n \\ & && -\xi \circ x \leq \varepsilon \mathbf{1}_n \\ &\text{and} && 0 \leq \xi \leq \mathbf{1}_n \end{aligned} \tag{36}$$

must remain bounded as $\varepsilon \downarrow 0$. Consequently, any limit point of such solutions must be a globally optimal solution of the original ℓ_0 -norm minimization problem (1).

Proposition 8 There exist positive scalars $\bar{\varepsilon}$ and σ such that for any $\varepsilon \in (0, \bar{\varepsilon}]$, we have $\|x\| \leq \sigma$ for any globally optimal solution (x, ξ) to (36).

Proof. For sufficiently small $\varepsilon > 0$, the problem (36) is equivalent to

$$\begin{aligned} & \underset{x}{\text{minimize}} && \sum_{j: x_j \neq 0} \left(1 - \frac{\varepsilon}{|x_j|} \right) \\ & \text{subject to} && Ax \geq b; \quad Cx = d, \end{aligned} \quad (37)$$

because one of the relaxed constraints is active if $x_j \neq 0$. For each index subset \mathcal{J} of $\{1, \dots, n\}$ with complement \mathcal{J}^c , consider the problem:

$$\begin{aligned} & \underset{x}{\text{minimize}} && \sum_{j \in \mathcal{J}: x_j \neq 0} \left(1 - \frac{\varepsilon}{x_j} \right) + \sum_{j \in \mathcal{J}^c: x_j \neq 0} \left(1 + \frac{\varepsilon}{x_j} \right) \\ & \text{subject to} && Ax \geq b; \quad Cx = d \\ & \text{and} && x_{\mathcal{J}} \geq 0, \quad x_{\mathcal{J}^c} \leq 0. \end{aligned} \quad (38)$$

The problem (37) is equivalent to the union of the problems (38) over all such index sets \mathcal{J} in the sense that if x is an optimal solution of the former problem, then x must be an optimal solution of the latter problem for some \mathcal{J} ; conversely, among the optimal solutions of (38) over all \mathcal{J} , the one that gives the smallest optimal objective value must be an optimal solution of (37). Therefore, it suffices to show that for each \mathcal{J} , the optimal solutions of (38) are uniformly bounded as $\varepsilon \downarrow 0$. Let $\mathcal{P}(\mathcal{J})$ denote the feasible region of (38) and let $\{\hat{x}^p\}_{p=1}^P$ and $\{\tilde{x}^r\}_{r=1}^R$ denote, respectively, the finite family of extreme points and rays of $\mathcal{P}(\mathcal{J})$. Any feasible solution of (38) can then be written as a convex combination of $\{\hat{x}^p\}_{p=1}^P$ plus a conic combination of $\{\tilde{x}^r\}_{r=1}^R$. In particular, if x^ε denotes an optimal solution of (38), then

$$x^\varepsilon = \sum_{p=1}^P \lambda_p^\varepsilon \hat{x}^p + \sum_{r=1}^R \mu_r^\varepsilon \tilde{x}^r$$

for some non-negative scalars $\{\lambda_p^\varepsilon\}_{p=1}^P$ and $\{\mu_r^\varepsilon\}_{r=1}^R$ with $\sum_{p=1}^P \lambda_p^\varepsilon = 1$. Clearly, $\hat{x}_{\mathcal{J}}^p \geq 0$, $\tilde{x}_{\mathcal{J}}^r \geq 0$, $\tilde{x}_{\mathcal{J}^c}^r \leq 0$, and $\tilde{x}_{\mathcal{J}^c}^r \leq 0$. Due to the sign consistency of the components, it follows that

$$|x^\varepsilon| = \sum_{p=1}^P \lambda_p^\varepsilon |\hat{x}^p| + \sum_{r=1}^R \mu_r^\varepsilon |\tilde{x}^r|.$$

Let

$$\rho \triangleq \min_{1 \leq r \leq R} \min_{1 \leq j \leq n} \{|\tilde{x}_j^r| : \tilde{x}_j^r \neq 0\}.$$

We claim that $\mu_r^\varepsilon \leq \varepsilon/\rho$ for all $\varepsilon > 0$ sufficiently small and all r . Indeed, if $\mu_{r_0}^\varepsilon > \varepsilon/\rho$ for some r_0 , then by defining

$$\tilde{x} \triangleq \sum_{p=1}^P \lambda_p^\varepsilon \hat{x}^p + \sum_{r_0 \neq r=1}^R \mu_r^\varepsilon \tilde{x}^r + \frac{\varepsilon}{\rho} \tilde{x}^{r_0},$$

we have

$$|x^\varepsilon| \geq |\tilde{x}|, \quad |x^\varepsilon| \neq |\tilde{x}|,$$

which implies

$$\sum_{j: x_j \neq 0} \left(1 - \frac{\varepsilon}{|x_j^\varepsilon|} \right) > \sum_{j: x_j \neq 0} \left(1 - \frac{\varepsilon}{|\tilde{x}_j|} \right),$$

contradicting the optimality of x^ε . □

5 Computational Results

After having discussed theoretical properties of the NLP reformulations, we now examine the practical performance of NLP solvers as solution methods of ℓ_0 -norm problems. One premise for our experiments is that black-box NLP codes are used with default settings. Those are applied directly to the NLP reformulations described in the previous sections, without modifications, despite the fact that some of these optimization models are not well-posed (i.e., the MFCQ does not hold at any feasible point for (17) and (18)).

The goal of these brute-force experiments is to assess the potential of NLP algorithms as solution approaches for hard ℓ_0 -norm optimization problems. If these initial experiments give encouraging results, it motivates further research that aims at a deeper understanding of the underlying mechanisms and the development of specialized methods.

The experiments were conducted using the NLP solvers CONOPT 3.15C [12], IPOPT 3.10.4 [28], KNITRO 8.0.0 [3], MINOS 5.51 [25], and SNOPT 7.2-8 [16]. We did not alter the solvers' default options, except that KNITRO was run with the option “`hessopt=5`”, which avoids the (potentially time-consuming) computation of the full Hessian matrix. In addition, any arising linear program (LP), mixed-integer linear programs (MILP), quadratic program (QP), and mixed-integer quadratic programs (MIQP) was solved with CPLEX 12.5.1.0. Matlab R2012b and the AMPL modeling software [15] were used as scripting languages and to generate the random problem instances. All numerical experiments reported in this paper were obtained on a 8-core 3.4GHz Intel Core i7 computer with 32GB of memory, running Ubuntu Linux.

5.1 Sparse solutions of linear inequalities

We first consider random instances of the model problem (1) of the form

$$\begin{aligned} & \underset{x \in \mathbb{R}^n}{\text{minimize}} && \|x\|_0 \\ & \text{subject to} && Ax \geq b \quad \text{and} \quad -M\mathbf{1}_n \leq x \leq M\mathbf{1}_n, \end{aligned} \tag{39}$$

where $A \in \mathbb{R}^{m \times n}$, $b \in \mathbb{R}^m$, and $M > 0$. The test instances were generated using AMPL's internal random number generator, where the elements of A and b are independent uniform random variables between -1 and 1.

Our numerical experiments compare the performance of different NLP optimization codes when they are applied to the different NLP reformulations. Because these problems are nonconvex, we also explore the effect of different starting points.

The NLP reformulations considered are:

- “Half”: The half-complementarity formulation (6);
- “Aggregate”: The full complementarity formulation where the complementarity constraints are reformulated by the inner product (17);
- “Individual”: The full complementarity formulation where the complementarity constraints are reformulated by the Hadamard product (18);
- “Squared”: The full complementarity formulation where the complementarity constraints are reformulated using the square of the inner product (20).

In addition to these, which were discussed in detail in the previous sections, we present results for the following formulations:

- “AMPL”: This formulation uses the keyword `complements` in order to pose (3) directly as an LPCC in AMPL. It is then up to the particular chosen optimization code to handle the complementarity constraints appropriately. Among the solvers considered here, only KNITRO is able to handle the `complements` keyword. KNITRO then reformulates the complementarity constraints using a penalty term that is added to the objective function; see [22] for details.
- “MILP”: The MILP formulation (7). Usually, a suitable choice of the upper bound M is not known. However, our test problems (39) explicitly include a bound on the optimal value x^* , so that the same M can be used in (39) and (7). The solution for this formulation is the global solution for (1).
- “LP”: The linear programming formulation (2). This is not an equivalent reformulation of the (1); instead, it is its commonly used convex approximation.

Because the nonlinear optimization methods aim at finding only local (and not global) optimal solutions of the nonconvex NLP reformulations, the choice of the starting point is crucial. In the experiments, we considered the following options:

- Start1: Set $x^+ = x^- = 0$ and $\xi = 0$;
- Start2: Set $x^+ = x^- = 0$ and $\xi = \mathbf{1}_n$;
- Start3: Let x^{LP} be the optimal solution of the LP formulation (2). Then set $x^+ \triangleq \max\{0, x^{\text{LP}}\}$, $x^- \triangleq \max\{0, -x^{\text{LP}}\}$, and $\xi = 0$;
- Start4: Let x^{LP} be the optimal solution of the LP formulation (2). Then set $x^+ \triangleq \max\{0, x^{\text{LP}}\}$, $x^- \triangleq \max\{0, -x^{\text{LP}}\}$, and $\xi = \mathbf{1}_n$;
- Start5: Let x^{LP} be the optimal solution of the LP formulation (2). Then set $x^+ \triangleq \max\{0, x^{\text{LP}}\}$, $x^- \triangleq \max\{0, -x^{\text{LP}}\}$, and ξ according to (4).

5.1.1 Pilot study on small problems

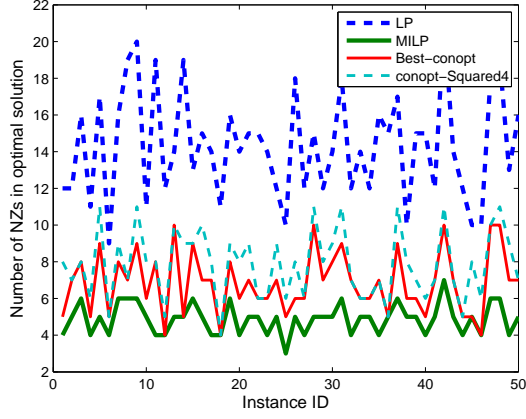
As a pilot study, we considered small problems with 30 constraints and 50 variables (i.e., $m = 30$, $n = 50$), and M was chosen to be 1000. To make statements with statistical significance, we generated 50 different random instances. Each of these instances is solved by 90 combinations of NLP solver, problem reformulation, and starting point, in addition to the LP and MILP formulations.

For each individual run, the point x^* returned by the solver is accepted as solution if it satisfies $Ax^* \geq b$, independent of the solvers’ exit status. The number of nonzeros (i.e., $\|x^*\|_0$) is computed by counting the number of elements with $|x_j^*| > 10^{-6}$. Table 1 lists the mean and standard deviation of $\|x^*\|_0$ for the different combinations. We note that all 50 problems were solved (i.e., the returned point satisfied the linear inequalities) for each combination, except for one CONOPT combination. As a reference, for the LP option, the mean was 14.32 with standard deviation 2.94, and for the MILP option, the mean was 4.88 with standard deviation 0.82.

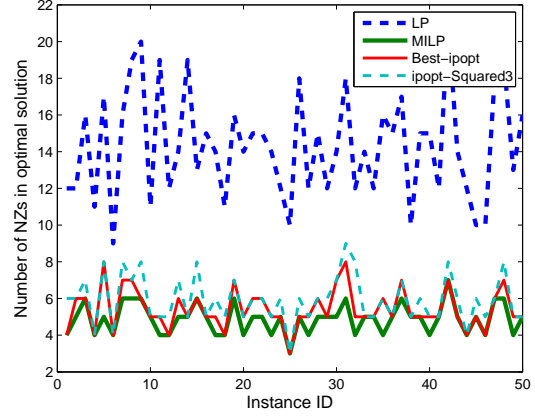
To present the results in more detail, Figures 1(a)–1(e) depict, for each of the 50 instances, the best ℓ_0 -norm obtained by the different NLP solvers, in comparison to the LP approximation and the (globally optimal) MILP solution. For each solver, the “Best-” line is the smallest $\|x^*\|_0$ value

		Start1		Start2		Start3		Start4		Start5	
		Mean	StdDev	Mean	StdDev	Mean	StdDev	Mean	StdDev	Mean	StdDev
CONOPT	Aggregate	22.32	5.35	N/A	all failed	14.32	2.94	8.98	2.17	14.32	2.94
	Individual	19.36	3.10	19.36	3.10	14.32	2.94	14.32	2.94	14.32	2.94
	Squared	10.34	3.42	8.82	2.14	9.16	2.85	7.86	1.89	8.72	2.62
IPOPT	Half	19.74	4.49	8.40	2.17	10.76	3.00	7.62	1.77	10.78	2.89
	Aggregate	6.32	1.41	6.78	1.65	6.26	1.35	6.50	1.57	6.50	1.59
	Individual	6.54	1.43	7.14	1.81	6.46	1.42	6.82	1.66	6.94	1.67
KNITRO	Squared	5.98	1.26	6.58	1.51	5.92	1.28	6.34	1.54	6.48	1.51
	Half	8.40	1.94	8.06	1.76	7.46	1.78	7.48	1.70	7.46	1.88
	Aggregate	10.08	3.77	9.30	3.21	16.50	2.97	15.82	3.30	15.78	2.60
MINOS	Individual	7.06	1.82	7.66	6.27	7.18	1.34	7.12	1.46	7.28	1.33
	Squared	6.30	1.35	6.40	1.55	8.14	1.93	8.56	2.37	7.78	1.97
	AMPL	9.52	7.35	9.78	5.88	9.44	4.41	9.12	4.26	9.04	4.60
SNOPT	Aggregate	15.08	2.59	10.02	2.49	14.32	2.94	14.32	2.94	14.32	2.94
	Individual	17.24	2.61	17.24	2.61	14.56	3.19	14.54	3.26	14.54	3.26
	Squared	9.16	2.21	7.86	1.84	9.04	2.90	7.50	1.81	7.44	1.81
SNOPT	Aggregate	11.04	2.46	9.18	2.25	14.32	2.94	14.26	2.87	14.32	2.94
	Individual	17.92	2.86	17.92	2.86	17.38	3.53	17.52	3.54	17.52	3.54
	Squared	7.94	1.78	7.60	1.71	7.76	2.15	8.00	2.27	7.56	1.68

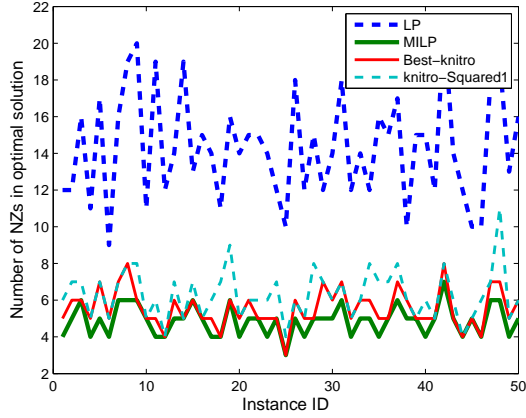
Table 1: Solution quality statistics for pilot study, grouped by solvers.



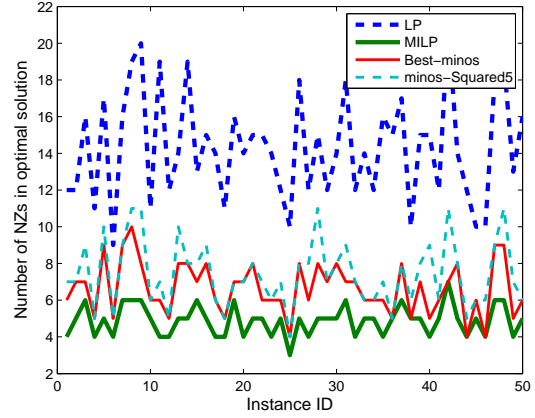
(a) CONOPT



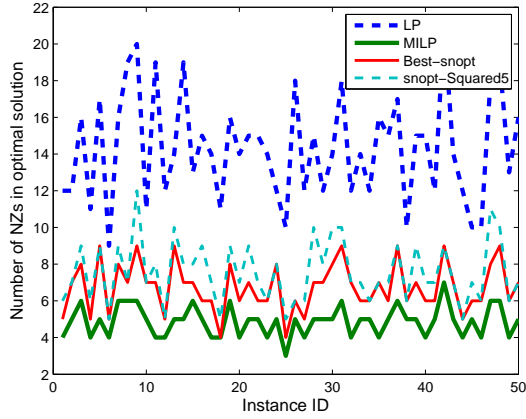
(b) IPOPT



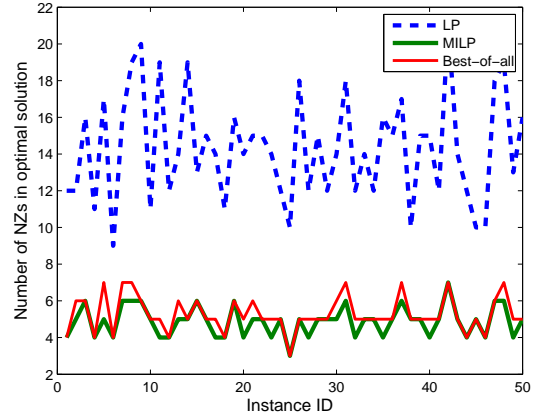
(c) KNITRO



(d) MINOS



(e) SNOPT



(f) Best of all solvers

Figure 1: The sparsest solutions for the 50 random problems using different NLP solvers in the pilot study.

obtained over *any* of the several combinations. In addition, the figures show the results obtained with the formulation/starting-point combination giving the smallest mean by the respective NLP solver. For example, in Figure 1(a), the line “conopt-Squared4” shows the outcome for the squared formulation and Start4 with the CONOPT solver. As we can see, all of the NLP solvers are able to find solutions that are sparser than those obtained by the common ℓ_1 -approximation. Indeed, the optimal solutions of some NLP solvers, particularly IPOPT and KNITRO, are able to find points with sparsity very close to the sparsest solution possible, as computed by the MILP formulation. Finally, Figure 1(f) shows the sparsest solution obtained by *any* of the solver, formulation, and starting point combinations. We see that, for each instance, at least one combination resulted in a solution that is equal or at most two nonzero elements worse than the global solution.

These results indicate that the application of (standard) NLP solvers to complementarity formulations of the ℓ_0 -problem results in high-quality solutions, considerably better than what is obtained by the common ℓ_1 -approximation. This promising observation is noteworthy, given that the problems are highly nonconvex and that there are many (undesirable) KKT points (see Corollary 1) to which the NLP solvers, which aim to find a KKT point, might potentially converge. The successful computational results suggest that the observations for the example in Section 4.1 may hold more broadly.

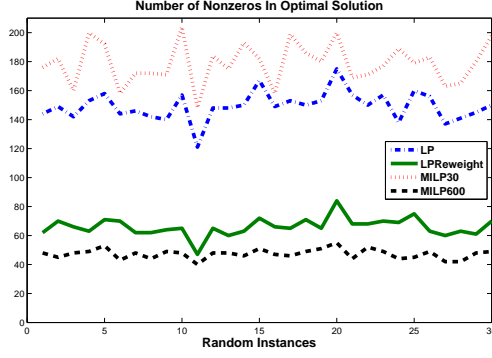
It is also somewhat surprising that, in this preliminary experiment, the squared formulation resulted in the lowest sparsity on average, despite the fact that no KKT point exists for any instance (see Proposition 4). The fact that the solvers terminate nevertheless can be explained by looking at the optimality conditions of (20). In particular, the reason why no KKT point exists is that the first condition in (21) cannot be satisfied exactly. However, as the NLP solver converges, the quantity on the right-hand-side of the complementarity in that condition can become arbitrarily small, as long as $(x^+, x^-) \neq 0$ and μ converges to infinity. So, even though there is no finite KKT point, the NLP solvers’ termination tests can be satisfied by diverging multipliers.

We also note that the CONOPT, MINOS, and SNOPT solvers usually do not converge to good solutions for the Aggregate and Individual formulations when started from a point obtained by the LP formulation (Start3–Start5). Indeed, in many cases the solvers terminate immediately at such a starting point. This can be explained by the fact that any feasible point is a KKT point for these formulations, and the active set solvers simply compute the corresponding multipliers at the starting point, so that the termination test is immediately satisfied. This is in contrast to the interior point solvers IPOPT and KNITRO, which are required to move the starting point away from bound constraints. This modification results in violations of the respective reformulation of the complementarity constraints and forces the algorithm to take steps.

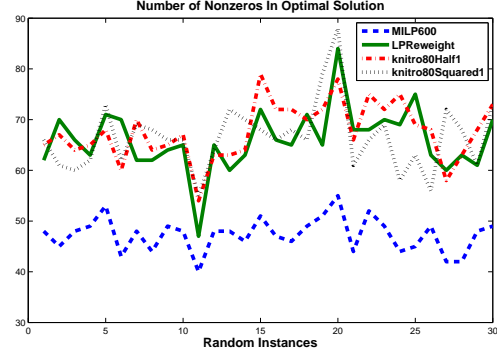
5.1.2 Large-scale problems

Based on the results in the pilot study, we pursued further numerical studies on larger problems, using the NLP solvers IPOPT and KNITRO. For this set of experiments we generated 30 random instances of (39) with 300 constraints and 1,000 variables.

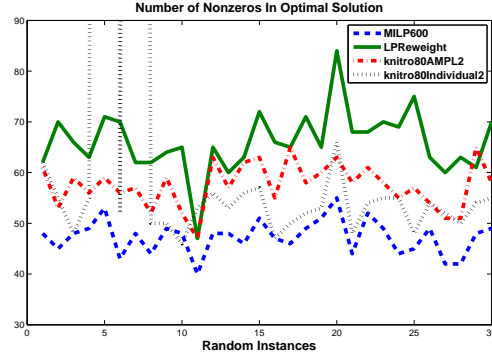
With this problem size, obtaining the true global optimum with the MILP formulation (7) is not possible with reasonable computational effort, even though we chose a rather small big-M constant ($M = 10$) in (39). In order to get an idea of what the sparsest solution for an instance might be, we ran CPLEX in multi-threaded mode for 10 minutes, which is equivalent to more than an hour of CPU time (this option is labeled MILP600), and we report the sparsity of the best incumbent.



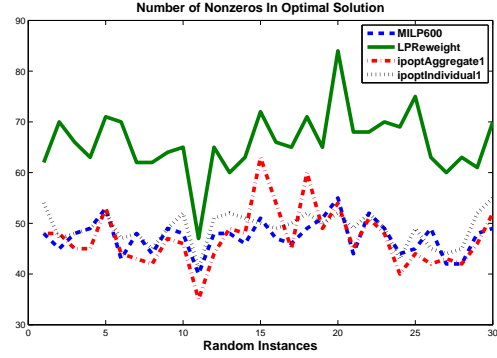
(a) LP/MILP



(b) KNITRO



(c) KNITRO



(d) IPOPT

Figure 2: The sparsest solutions for the 30 random problems using different NLP solvers in the large-scale study.

Similarly, we also wanted to explore the quality of a heuristic solution that an MILP solver is able to find in a time that is comparable to that taken by the NLP solvers. For this purpose, the MILP60 option reports the best incumbent obtained in one minute, equivalent to about 2.5 minutes of CPU time. In these experiments, CPLEX was run with the `mipemphasis=1` option, to focus on finding good heuristic solutions quickly.

As we observed in the small-case study, the ℓ_1 -approximation (2) did not lead to good solutions. However, it is common to enhance the LP solution by some improvement heuristic. One such approach is the iterative re-weighted ℓ_1 -minimization scheme proposed in [6]. Starting with the optimal solution $x^{*,0}$ of (2), this procedure optimizes a sequence of LPs for $k = 0, 1, 2, \dots$ to generate iterates from

$$x^{*,k+1} = \operatorname{argmin}_{x \in \mathbb{R}^n} \left\{ \sum_{i=1} w_{k,i} |x_i|, \quad \text{subject to } Ax \geq b \text{ and } -M\mathbf{1}_n \leq x \leq M\mathbf{1}_n \right\},$$

with weights $w_{k,i} = 1/|x_i^{*,k}|$. Here, we understand $w_{k,i} = 1/0 = \infty$ as x_i being fixed to zero. We ran this procedure for 30 iterations (after which the iterates had settled), and report the outcome of this procedure under the label “LPReweight”.

The results of the experiments are depicted in Figure 2. First, we see in Figure 2(a) that the iterative re-weighting procedure indeed improves the straight-forward ℓ_1 -approximation considerably; it more

than halves the objective function. However, there is still a significant gap (17% – 66%) between LPReweight and the best solution found by the MILP solver within an hour of CPU time. We note that the MILP solver is not able to find any good solution within about 2.5 minutes of CPU time.

Figures 2(b)–2(c) show the solution quality obtained with different reformulations of the ℓ_0 -problem when solved with KNITRO. For each formulation, we report the best outcome among Start1 and Start2; to limit the amount of data in the graphs, we plot only selected representative combinations, including the best ones. In this experiment, the Squared and Half formulation are able to generate solutions of similar sparsity to those obtained by the LPReweight option. The AMPL and Individual formulations give better solutions, with the latter closing the gap considerably, except for two failures for instances 5 and 7. IPOPT results are reported in Figure 2(d) for the Aggregate and Individual formulations. We see that the solution quality is comparable with that obtained by the MILP solver, and for the Aggregate formulation often even better.

For practical purposes it is also important to consider the computation time required to solve the NLP formulations. Table 2 lists the average solution quality and required CPU time for representative combinations of formulations and starting points. We note that, on average, solutions within 4% of the MILP600 objective can be computed in less than one minute (“ipoptIndividual1”), and solutions comparable and better than the MILP600 objective can be computed in less than 10 minutes (“ipoptAggregate1”).

	$\ x^*\ _0$		CPU time	
	Mean	StdDev	Mean	StdDev
LP	149.33	9.93	6.56	2.29
LPReweight	66.00	6.22	≈ 6.56	≈ 2.29
MILP30	178.87	14.21	155.90	15.11
MILP600	47.33	3.41	4011.55	302.77
knitro80Individual2*	53.07	4.33	586.04	824.28
knitro80Aggregate1	76.03	4.21	89.16	59.00
knitro80AMPL2	57.50	4.42	48.59	7.80
knitro80Half1	67.63	5.66	23.91	8.53
knitro80Squared1	66.53	6.71	326.77	129.06
ipoptIndividual1	49.07	3.31	53.17	23.90
ipoptAggregate1	47.17	5.74	577.62	284.24

Table 2: Summary statistics for large-scale study. (* results reported with the two outliers removed.)

5.2 Tracing a Pareto curve

In this section, we are concerned with the Pareto curve defined by the problem

$$\begin{aligned}
& \underset{x \in \mathbb{R}^n}{\text{minimize}} && q(x) + \gamma \|x\|_0 \\
& \text{subject to} && Ax \geq b
\end{aligned} \tag{40}$$

with varying the weighting parameter $\gamma > 0$, where $q(x)$ is a convex function, $A \in \mathbb{R}^{n \times m}$, and $b \in \mathbb{R}^m$. Similar to the LP approximation (2) of (1), the discontinuous ℓ_0 -norm in (41) may be

approximated by the convex ℓ_1 -norm to give

$$\begin{aligned} & \underset{x \in \mathbb{R}^n}{\text{minimize}} && q(x) + \gamma \|x\|_1 \\ & \text{subject to} && Ax \geq b. \end{aligned} \tag{41}$$

Variations of this convex approximation are found frequently in the literature, for example in the context of compressive sensing and basis pursuit denoising [2] and in image deblurring [1].

In this section, we explore how well the Pareto curve is traced if the points are computed by NLP solvers that are applied to the NLP reformulations of the complementarity formulation of (40):

$$\begin{aligned} & \underset{x, x^\pm, \xi}{\text{minimize}} && q(x) + \gamma \mathbf{1}_n^T (\mathbf{1}_n - \xi) \\ & \text{subject to} && Ax \geq b, \quad x = x^+ - x^- \\ & && 0 \leq \xi \perp x^+ + x^- \geq 0 \\ & && 0 \leq x^+ \perp x^- \geq 0 \\ & \text{and} && \xi \leq \mathbf{1}_n. \end{aligned} \tag{42}$$

In our experiment, we considered a random instance with 200 constraints and 500 variables. Here, $q(x) = \frac{1}{2}x^T Qx + c^T x$ with $Q = CC^T$ is a strictly convex quadratic function, and all the elements of A , b , C , and c were chosen uniformly in $[-1, 1]$ at random.

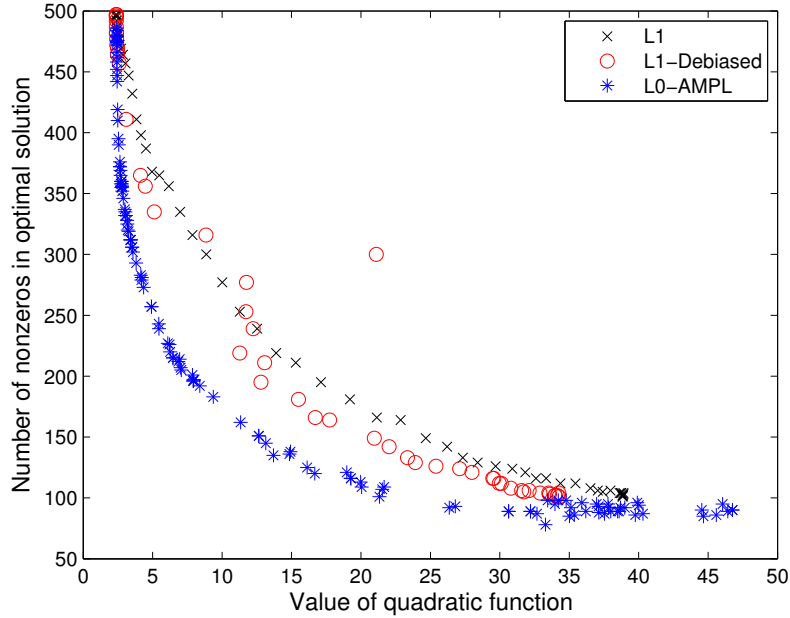


Figure 3: Pareto curve for a random instance with 200 constraints and 500 variables.

Figure 3 depicts the points obtained by three different methods in the $(q(x), \|x\|_0)$ space. The curve labeled “L1” is obtained by solving the QP formulation of (41) with γ varying in the interval $[10^{-2}, 10^4]$. The “L1-Debiased” points are obtained by a standard “debiasing” technique for improving the “L1” solutions (see, e.g., [29]). This procedure takes the optimal solution x^* of (41),

determines the set of fixed variables $\mathcal{J} = \{j : |x_j^*| \leq \epsilon\}$ (we use the tolerance $\epsilon = 10^{-7}$) and resolves the original problem (41) without a regularization term ($\gamma = 0$) and all variables with indices in \mathcal{J} fixed to zero. In other words, the ℓ_1 -regularized problem (41) is solved first to determine the sparsity pattern for the variables, and then the “L1-Debiased” solution is the best objective value that can be achieved with that pattern. By this definition, the points for the “L1-Debiased” solution should always be to the left of the “L1” curve. The exceptions to this rule are due to convergence failures of CPLEX¹.

The third set of points (“L0-AMPL”) is obtained by solving the AMPL formulation of the ℓ_0 -regularized problem (42) with KNITRO using the Start2 starting point. Because the termination test in KNITRO 8.0.0 led to premature termination in some cases, we ran this experiment with the options “hessopt=5 maxit=3000 feastol=0 feastol_abs=1e-6 opttol=0 opttol_abs=1e-6”. The computation times for the KNITRO runs vary between 37–167 CPU seconds. The weighting parameter γ was varied between $[0.01, 10]$, with a high concentration of values within the interval $[0.1, 0.5]$ where most of the changes occurred.

In this experiment, the NLP solver achieves much sparser solutions than the LP-based formulations for a given value of the quadratic function $q(x)$. We also observe that the ℓ_1 -formulation cannot produce solutions with fewer than 101 nonzeros, while the NLP-based approach is able to generate solutions with as few as 78 nonzero elements.

We point out that the set of local minimizer of the ℓ_0 -norm problem (40) is independent of the choice of γ (see Proposition 2). Therefore, it is not clear that an NLP solver, which aims to find local solutions, is able to trace a Pareto-type curve as γ is varied. However, as we see in Figure 3, it appears that different values of the weighting parameter guide the NLP solver to different local solutions that turn out to trace the Pareto curve quite well. In the experiment, the number of nonzeros in the solution tends to decrease with increasing values of γ ; however, this relationship is not strictly monotone.

We also stress that the outcome of the experiment is rather sensitive to the choice of the starting point, NLP solver, and the reformulation of the ℓ_0 -norm. In our tests, most configurations were not able to generate a clear Pareto curve, and we report here only the most successful one.

6 Conclusions and Outlook

We presented several nonlinear programming reformulations of the ℓ_0 -norm minimization problem. Our numerical study suggests that standard NLP codes are able to generate solutions that can be significantly better than those obtained by convex approximations of the NP-hard problem. This is somewhat remarkable because the NLP formulations are highly nonconvex and the usual constraint qualifications, such as MFCQ, do not hold.

Typically, NLP algorithms are designed to find a KKT point, ideally one that satisfies the second-order necessary optimality conditions. Our analysis pertaining to the optimality conditions of the NLP formulations finds that, for the simple problem with linear constraints in the introduction, any feasible point for the ℓ_0 -norm minimization problem is such a KKT point. Consequently, from this perspective, any feasible point seems equally attractive to the NLP solver, and therefore these considerations do not explain the observed high quality of the solutions.

¹In these cases, CPLEX’s exit message was: “best solution found, primal-dual infeasible” or similar.

We also discussed the properties of solutions for relaxations of the NLP formulations as the relaxation parameter is driven to zero. We established that global minimizers of the relaxed problem converge to a global minimizer of original problem. This result justifies the choice of the relaxation. Furthermore, for a small example problem with linear constraints we showed that there are only a few KKT points for the relaxed problem, and that those converge to a small number of limit points as the relaxation parameter goes to zero. This is in contrast to the earlier result that does not distinguish between any two feasible points. In addition, we established that a KKT point for the relaxed NLP that is not close to a global minimizer of the original problem is a local *maximizer* for the relaxation. As a consequence, an NLP solver, when applied to the relaxation with a small relaxation parameter, will most likely converge to a point that is close to a global minimizer of the original ℓ_0 -norm problem.

These observations might help to explain why the NLP solvers compute points with objective values close to the globally optimal value in our experiments. Some solvers relax any given NLP by a small amount by default, and therefore explicitly solve a relaxation of the complementarity reformulation. For other solvers, numerical inaccuracies or the linearization of the nonlinear reformulations of the complementarity constraints at infeasible points might have an effect that is similar to that of a relaxation. The details of such an analogy, as well as the generalization of the results beyond the particular small example, are subject to future research.

Our numerical experiments did not identify a clear winner among the different reformulations of the ℓ_0 -norm minimization problems. Similarly, while some NLP codes tended to produce better results than others, it is not clear which specific features of the algorithms or their implementations are responsible for finding good solutions. We point out that each software implementation includes enhancements, such as tricks to handle numerical problems due to round-off error or heuristics that are often not included in the mathematical description in scientific papers. Because the NLP reformulations of the ℓ_0 -problems are somewhat ill-posed, these enhancement are likely to be crucial for the solver's performance. Once the relevant ingredients of the reformulation and optimization method have been identified, it might be possible to design specialized NLP-based algorithms that are tailored to the task of finding sparse solutions efficiently.

Finally, the numerical study in this paper has been performed using randomly generated model problems. Future efforts will explore the suitability of the proposed approach for ℓ_0 -norm optimization problems arising in particular applications areas including compressive sensing [5, 9], basis pursuit [2, 8], machine learning [7, 24], and genome-wide association studies [30].

References

- [1] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, 2(1):183–202, 2009.
- [2] E. van den Berg and M. P. Friedlander. Probing the Pareto frontier for basis pursuit solutions. *SIAM Journal on Scientific Computing*, 31(2):890–912, 2008.
- [3] R. H. Byrd, J. Nocedal, and R. A. Waltz. KNITRO: An Integrated Package for Nonlinear Optimization. In *Large-Scale Nonlinear Optimization*, edited by G. di Pillo and M. Roma. Springer-Verlag, pp. 35–59, 2006.

- [4] E. Candès, M. Rudelson, T. Tao, and R. Vershynin. Error correction via linear programming. In *46th Annual IEEE Symposium on Foundations of Computer Science, 2005. FOCS 2005*, pages 668–681, 2005.
- [5] E. J. Candès and M. Wakin. An introduction to compressive sampling. *IEEE Signal Processing Magazine*, 25(2):21–30, 2008.
- [6] E. J. Candès, M. Wakin, and S. Boyd. Enhancing sparsity by reweighted l_1 minimization. *Journal of Fourier Analysis and Applications*, 14:877–905, 2007.
- [7] C. Chen and O. L. Mangasarian. Hybrid misclassification minimization. *Advances in Computational Mathematics*, 5:127–136, 1996.
- [8] S. S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. *SIAM Journal on Scientific Computing*, 20(1):33–61, 1998.
- [9] M. Davenport, M. Duarte, Y. Eldar, and G. Kutyniok. Introduction to compressed sensing. Chapter 1 of *Compressed Sensing: Theory and Applications*, edited by Y. Eldar, and G. Kutyniok. Cambridge University Press, 2012.
- [10] A.-V. de Miguel, M. Friedlander, F. J. Nogales, and S. Scholtes. A two-sided relaxation scheme for mathematical programs with equilibrium constraints. *SIAM Journal on Optimization*, 16(2):587–609, 2006.
- [11] D. L. Donoho and M. Elad. Optimally sparse representation in general (nonorthogonal) dictionaries via ℓ_1 minimization. *Proceedings of the National Academy of Sciences*, 100(5):2197–2202, 2003.
- [12] A. Drud. CONOPT: A GRG code for large sparse dynamic nonlinear optimization problems. *Mathematical Programming*, 31(2):153–191, 1985.
- [13] F. Facchinei and J. S. Pang. *Finite-dimensional variational inequalities and complementarity problems: Volumes I and II*. Springer-Verlag, New York, 2003.
- [14] R. Fletcher and S. Leyffer. Solving mathematical programs with complementarity constraints as nonlinear programs. *Optimization Methods and Software*, 19(1):15–40, 2004.
- [15] R. Fourer, D. M. Gay, and B. W. Kernighan. *AMPL: A modeling language for mathematical programming*. Second edition. Thomson, Toronto, Canada, 2003.
- [16] P. E. Gill, W. Murray, and M. A. Saunders. SNOPT: An SQP algorithm for large-scale constrained optimization. *SIAM Review*, 47(1):99–131, 2005.
- [17] R. Gribonval and M. Nielsen. Sparse representations in unions of bases. *IEEE Transactions on Information Theory*, 49(12):3320–3325, 2003.
- [18] J. Hu, J. E. Mitchell, J. S. Pang, K. P. Bennett, and G. Kunapuli. On the global solution of linear programs with linear complementarity constraints. *SIAM Journal on Optimization*, 19(1):445–471, 2008.
- [19] J. Hu, J. E. Mitchell, J. S. Pang, and B. Yu. On linear programs with linear complementarity constraints. *Journal of Global Optimization*, 53(1):29–51, 2012.

- [20] R. Janin. Directional derivative of the marginal function in nonlinear programming. *Mathematical Programming Study* 21 (1984) 110126.
- [21] C. Kanzow and A. Schwartz. The price of inexactness: Convergence properties of relaxation methods for mathematical programs with equilibrium constraints revisited. Technical report, Institute of Mathematics, University of Würzburg, Würzburg, Germany, March 2013.
- [22] S. Leyffer, G. López-Calva, and J. Nocedal. Interior methods for mathematical programs with complementarity constraints. *SIAM Journal on Optimization*, 17(1):52–77, 2006.
- [23] Z.-Q. Luo, J. S. Pang, and D. Ralph. *Mathematical Programs with Equilibrium Constraints*. Cambridge University Press, Cambridge, 1996.
- [24] O. Mangasarian. Misclassification minimization. *Journal of Global Optimization*, 5:309–323, 1994.
- [25] B. A. Murtagh and M. A. Saunders. MINOS 5.5 User’s Guide. Report SOL 83–20R, Systems Optimization Laboratory, Stanford University, December 1983 (revised February 1995).
- [26] S. Scholtes. Convergence properties of a regularisation scheme for mathematical programs with complementarity constraints. *SIAM Journal on Optimization*, 11(4):918–936, 2001.
- [27] R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society B*, 58(1):267–288, 1996.
- [28] A. Wächter and L. T. Biegler. On the implementation of a primal-dual interior point filter line search algorithm for large-scale nonlinear programming. *Mathematical Programming*, 106(1):25–57, 2006.
- [29] S. J. Wright, R. D. Nowak, and M. A. T. Figueiredo. Sparse reconstruction by separable approximation. *IEEE Transactions on Signal Processing*, 57(7):2479–2493, 2009.
- [30] T. T. Wu, Y. F. Chen, T. Hastie, E. Sobel, and K. Lange. Genome-wide association analysis by lasso penalized logistic regression. *Bioinformatics*, 25(6):714–721, 2009.
- [31] Y. Zhang. Theory of compressive sensing via ℓ_1 -minimization: a non-RIP analysis and extensions. *Journal of the Operations Research Society of China*, 1(1):79–105, 2013.