

Nonsmooth Newton's Method: Some Structure Exploitation

Alberto De Marchi^[0000-0002-3545-6898] and Matthias
Gerdts^[0000-0001-8674-5764]

Department of Aerospace Engineering
Bundeswehr University Munich
Werner-Heisenberg-Weg 39
85577 Neubiberg, Germany
alberto.demarchi@unibw.de
matthias.gerdts@unibw.de

Abstract. We investigate real asymmetric linear systems arising in the search direction generation in a nonsmooth Newton's method. This applies to constrained optimisation problems via reformulation of the necessary conditions into an equivalent nonlinear and nonsmooth system of equations. We propose a strategy to exploit the problem structure. First, based on the sub-blocks of the original matrix, some variables are selected and ruled out for a posteriori recovering; then, a smaller and symmetric linear system is generated; eventually, from the solution of the latter, the remaining variables are obtained. We prove the method is applicable if the original linear system is well-posed. We propose and discuss different selection strategies. Finally, numerical examples are presented to compare this method with the direct approach without exploitation, for full and sparse matrices, in a wide range of problem size.

Keywords: Structure Exploitation · Linear Algebra · Nonsmooth Newton's Method · Nonlinear Optimization.

1 Introduction

In this paper, we consider the real square nonsymmetric possibly large sparse linear system

$$\begin{bmatrix} \mathbf{Q} & \mathbf{A}^\top & \mathbf{C}^\top \\ \mathbf{A} & & \\ -\mathbf{S}\mathbf{C} & & \mathbf{T} \end{bmatrix} \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \\ \mathbf{z} \end{pmatrix} = \begin{pmatrix} \mathbf{f} \\ \mathbf{g} \\ \mathbf{h} \end{pmatrix} \quad (1)$$

where $\mathbf{Q} \in \mathbb{R}^{n_x \times n_x}$, $\mathbf{A} \in \mathbb{R}^{n_a \times n_x}$, $\mathbf{C} \in \mathbb{R}^{n_c \times n_x}$ and $\mathbf{S}, \mathbf{T} \in \mathbb{R}^{n_c \times n_c}$ are given matrices and $\mathbf{f} \in \mathbb{R}^{n_x}$, $\mathbf{g} \in \mathbb{R}^{n_a}$, $\mathbf{h} \in \mathbb{R}^{n_c}$ are given vectors (n_x, n_a, n_c being some positive integers); \mathbf{S} and \mathbf{T} are non-zero diagonal matrices. The contribution of this paper is the exploitation of the structure in problem (1) and its transformation into a smaller symmetric linear system, with a saddle-point structure, from which the solution to (1) can be easily recovered. In particular, two stages are discussed. First, a reduction step generates a smaller linear system and a way

to recover eliminated variables from the solution of this reduced system. This step exploits the fact that matrix \mathbf{T} is diagonal, and then symbolically solve for (some of) the variables in \mathbf{z} . Several different reduction strategies are discussed and compared. The second step aims at rewriting the linear system in a symmetric form, allowing to adopt solvers for symmetric systems, which are usually more time and memory efficient. Despite these advantages, some computational overhead is needed, especially in the reduction step, which might introduce a break-even point, that is, this exploitation may pay off, e.g., in terms of computational time, only under certain conditions, e.g., large instances. In fact, an optimal reduction strategy might exist, depending on the specific properties of the problem; indeed, it may even depend on the specific values of the entries. Throughout the paper, we investigate the influence of the reduction strategy on the performance of the aforementioned two-steps exploitation; however, a detailed optimization of the reduction policy is beyond the scope of this paper. We point out that the proposed method could be combined with constraint-reduction approaches as, e.g., those presented in [5,12].

Once the original problem (1), say $\mathbf{V}\mathbf{d} = \mathbf{r}$ for brevity, has been transformed, a reduced symmetric linear system, say $\hat{\mathbf{V}}\hat{\mathbf{d}} = \hat{\mathbf{r}}$, is to be solved. To this end, any method can be adopted. The choice may depend on the problem, in particular on its size, fill-in, sparsity pattern, accuracy requirements and memory constraints, availability of good preconditioners, and so on. Within this work, we compare the effectiveness and the limitations of the proposed method for structure exploitation when a direct solver is adopted to tackle the linear system.

1.1 Motivation

Linear systems with the form (1) arise, e.g., from nonlinear complementarity problems [6], nonlinear optimization problems with inequality constraints [7] and discretized optimal control problems with state and control constraints [8,9]. Usually, these are reformulated through the Karush-Kuhn-Tucker (KKT) necessary optimality conditions, then equivalently transformed into a nonlinear system of equations with the so called NCP-functions [16] and finally solved with a nonsmooth version of Newton's method [14]. Some globalization strategies [9,11]

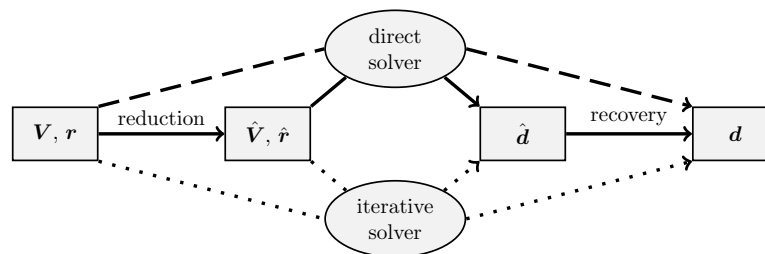


Fig. 1. Solution diagram: direct methods with (solid) and without (dashed) structure exploitation are compared; iterative methods (dotted) are not considered here.

and results in functions spaces [18] have been reported. It has been shown that, for some NCP-functions, this approach is equivalent to a primal-dual active set strategy [10]. Indeed, different NCP-functions exhibit different properties and might affect the convergence behaviour [1,16]. With reference to the problem (1), matrices \mathbf{Q} , \mathbf{A} and \mathbf{C} can be considered as iterate-dependent linear-quadratic approximations of an underlying nonlinear problem, while diagonal matrices \mathbf{S} and \mathbf{T} originate from the NCP-function adopted and vectors \mathbf{f} , \mathbf{g} and \mathbf{h} are the residuals of the aforementioned nonlinear system of equations.

Example Let us consider a quadratic program (QP) with linear equality and inequality constraints. Hence, we seek an $\mathbf{x} \in \mathbb{R}^{n_x}$ minimizing $\frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x} + \mathbf{q}^\top \mathbf{x}$, subject to constraints $\mathbf{A}\mathbf{x} = \mathbf{a}$, $\mathbf{C}\mathbf{x} \leq \mathbf{c}$, where \mathbf{Q} , \mathbf{A} , \mathbf{C} and \mathbf{q} , \mathbf{a} , \mathbf{c} are given matrices and vectors, respectively. Here the inequalities are understood componentwise. Linear constraints ensure that regularity conditions are met, then the KKT conditions are necessary for optimality; these read:

$$\mathbf{Q}\mathbf{x} + \mathbf{q} + \mathbf{A}^\top \boldsymbol{\lambda} + \mathbf{C}^\top \boldsymbol{\mu} = \mathbf{0} \quad (2a)$$

$$\mathbf{A}\mathbf{x} = \mathbf{a} \quad (2b)$$

$$\mathbf{0} \leq \boldsymbol{\mu} \perp \mathbf{C}\mathbf{x} \leq \mathbf{c} \quad (2c)$$

where $\boldsymbol{\lambda}$ and $\boldsymbol{\mu}$ denote the Lagrange multipliers. In (2c), inequality and complementarity constraints hold componentwise. Let us consider an NCP-function $\varphi : \mathbb{R}^2 \rightarrow \mathbb{R}$, e.g., the original or the penalized Fischer-Burmeister function [1,7], which by definition satisfies

$$\varphi(a, b) = 0 \Leftrightarrow 0 \leq a \perp b \geq 0 \quad (3)$$

for any pair (a, b) . Thanks to this property, the KKT system (2) is equivalently rewritten as a nonlinear system $\boldsymbol{\psi}(\mathbf{z}) = \mathbf{0}$, collecting vector $\mathbf{z} = (\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \in \mathbb{R}^{n_z}$, $n_z := n_x + n_a + n_c$, and with vector-valued function $\boldsymbol{\psi} : \mathbb{R}^{n_z} \rightarrow \mathbb{R}^{n_z}$ defined by

$$\boldsymbol{\psi}(\mathbf{z}) = \begin{pmatrix} \mathbf{Q}\mathbf{x} + \mathbf{q} + \mathbf{A}^\top \boldsymbol{\lambda} + \mathbf{C}^\top \boldsymbol{\mu} \\ \mathbf{A}\mathbf{x} - \mathbf{a} \\ \varphi(\mathbf{c} - \mathbf{C}\mathbf{x}, \boldsymbol{\mu}) \end{pmatrix} \quad (4)$$

Here the NCP-function φ applies componentwise. A (globalized) possibly nonsmooth Newton's-type method generates a sequence $\{\mathbf{z}^k\}$ through the recurrence $\mathbf{z}^{k+1} = \mathbf{z}^k + \alpha^k \mathbf{d}^k$, $k = 0, 1, 2, \dots$, where the step length $\alpha^k > 0$ is determined, e.g., by a line-search procedure of Armijo's type and the search direction \mathbf{d}^k is the solution of the linear equation $\mathbf{V}^k \mathbf{d} = -\boldsymbol{\psi}(\mathbf{z}^k)$ [7,8,9,11,14]. The matrix \mathbf{V}^k is an element of the Clarke's generalized Jacobian (that is the convex hull of the Bouligand differential [2]) of $\boldsymbol{\psi}$ at \mathbf{z}^k , namely $\mathbf{V}^k \in \partial \boldsymbol{\psi}(\mathbf{z}^k)$ [8,11]. The NCP-function φ is the only element in (4) which can possibly make function $\boldsymbol{\psi}$ nonsmooth. Hence, from (4), one obtains

$$\mathbf{V}^k = \begin{bmatrix} \mathbf{Q} & \mathbf{A}^\top & \mathbf{C}^\top \\ \mathbf{A} & & \\ -\mathbf{S}^k \mathbf{C} & & \mathbf{T}^k \end{bmatrix} \quad (5)$$

with diagonal matrices $\mathbf{S}^k = \text{diag}(s_1^k, s_2^k, \dots, s_{n_c}^k)$ and $\mathbf{T}^k = \text{diag}(t_1^k, t_2^k, \dots, t_{n_c}^k)$, whose entries are pairwise coupled via the (possibly generalized) differential of the NCP-function φ . Defining $\mathbf{v}^k := \mathbf{c} - \mathbf{C}\mathbf{x}^k$ the inequality constraint violation at the k -th iterate, for $i = 1, 2, \dots, n_c$, the coupling for the i -th inequality constraint reads [1,7,8]:

$$(s_i^k, t_i^k) \in \partial\varphi(v_i^k, \mu_i^k) \quad (6)$$

We remark that matrices \mathbf{S}^k and \mathbf{T}^k are diagonal because the NCP-function φ applies componentwise in (4). Then, the linear system $\mathbf{V}^k \mathbf{d} = -\boldsymbol{\psi}(\mathbf{z}^k)$ to compute the search direction \mathbf{d}^k corresponds exactly to problem (1).

1.2 Outline

This work is organized as follows. Section 2 outlines the structure exploitation strategy and introduces an underlying assumption. In Sections 2.1 and 2.2 the two main steps are developed and discussed. Section 3 validates the proposed approach numerically, for both full and sparse matrices, with different reduction strategies, showing effectiveness and limitations of the proposed algorithm. Section 4 concludes the paper and presents ideas for future research.

2 Structure exploitation

Problem (1), also denoted $\mathbf{V}\mathbf{d} = \mathbf{r}$ for brevity, can be directly solved via any linear algebra package, e.g., MA48 [4], PARDISO [15], SUPERLU [13], `mldivide` in MATLAB [17]. However, we aim at exploiting our knowledge about the structure of matrix \mathbf{V} ; computational effort and achieved accuracy might benefit from this, especially for large-scale and sparse linear systems. Firstly, we notice that matrix \mathbf{V} and vector \mathbf{r} are often computed blockwise and then assembled. Hence, matrices \mathbf{Q} , \mathbf{A} , \mathbf{C} and vectors \mathbf{s} , \mathbf{t} , \mathbf{f} , \mathbf{g} and \mathbf{h} are here considered as the starting ingredients for solving (1). Overall, two directions are explored, mainly exploiting the diagonal structure of \mathbf{S} and \mathbf{T} . In Section 2.1 a reduction step is discussed, eliminating some variables and introducing a smaller asymmetric linear system. Then, in Section 2.2, this is transformed into a symmetric one which is equivalent. Nonetheless, these operations for reorganizing the linear system and successive recovering of variables constitutes an overhead of computation. This means that these procedures might be worthy only for certain problems, likely large instances with lots of inequality constraints. In Section 3, numerical tests show that this break-even point corresponds to relatively small problem instances (for the tested implementation).

We point out that the proposed method relies on the following assumption on the diagonals of \mathbf{S} and \mathbf{T} ; it reads:

Assumption 1 For $i = 1, 2, \dots, n_c$, it holds $(s_i, t_i) \neq (0, 0)$.

This is a mild requirement, in that it corresponds to problem (1) to be well-posed. One can show the following result:

Lemma 1. *If problem (1) admits a unique solution, then Assumption 1 holds.*

Proof (by contradiction). Let us assume there exist \mathbf{d} unique solution to (1) and i such that $(s_i, t_i) = (0, 0)$. Hence, the row of \mathbf{V} corresponding to the i -th inequality constraint consists of zeros only. Then, the matrix \mathbf{V} is rank deficient and the problem is undetermined. Two cases are possible, depending on the value of h_i on the right-hand side. If $h_i = 0$, then (1) admits infinitely many solutions, hence solution \mathbf{d} is not unique. If $h_i \neq 0$, then the linear system is unsolvable (impossible) and \mathbf{d} cannot be a solution. \square

Remark 1. Assumption 1 requires a mild condition to be satisfied by the NCP-function φ . For instance, a sufficient condition is that for any given pair $(a, b) \in \mathbb{R}^2$ there exists a pair $(s, t) \in \partial\varphi(a, b)$ such that $(s, t) \neq (0, 0)$; this allows to choose always suitable entries for \mathbf{S} and \mathbf{T} . The Fischer-Burmeister function and the max function, among other NCP-functions, have this property.

Let us denote $\mathcal{I} := \{1, 2, \dots, n_c\}$ the index set for the inequality constraints, $\mathcal{I}_{\text{re}}^0 := \{i \in \mathcal{I} \mid t_i \neq 0\}$ and $\mathcal{I}_{\text{sy}}^0 := \{i \in \mathcal{I} \mid s_i \neq 0\}$ the largest index sets that allow respectively the reduction step and the symmetrization step, discussed below. Thanks to Assumption 1, these satisfy $\mathcal{I}_{\text{re}}^0 \cup \mathcal{I}_{\text{sy}}^0 = \mathcal{I}$. Let us consider an index subset $\mathcal{I}_{\text{re}} \subseteq \mathcal{I}_{\text{re}}^0$, sufficiently large to satisfy $\mathcal{I}_{\text{re}} \cup \mathcal{I}_{\text{sy}}^0 = \mathcal{I}$. Then, for the associated complement $\mathcal{I}_{\overline{\text{re}}} := \mathcal{I} \setminus \mathcal{I}_{\text{re}}$, it holds $\mathcal{I}_{\overline{\text{re}}} \subseteq \mathcal{I}_{\text{sy}}^0$. With this construction, it is possible to apply the reduction step, ruling out a given set \mathcal{I}_{re} of variables, and subsequently the symmetrization step on the linear system with the remaining variables, namely those in $\mathcal{I}_{\overline{\text{re}}}$.

Remark 2. We stress that in general it is $\mathcal{I}_{\text{re}}^0 \cap \mathcal{I}_{\text{sy}}^0 \neq \emptyset$ and hence the choice of \mathcal{I}_{re} is not unique. This suggests there could be an optimal reduction strategy, possibly dependent on \mathbf{V} and with some degree of computation awareness. However, this issue is beyond the scope of this paper.

In Section 3, we compare the following definitions of \mathcal{I}_{re} through numerical investigations:

$$\mathcal{I}_{\text{re}}^t := \left\{ i \in \mathcal{I} \mid |t_i| \geq \epsilon \right\} \quad (7a)$$

$$\mathcal{I}_{\text{re}}^s := \left\{ i \in \mathcal{I} \mid |s_i| \leq \epsilon \right\} \quad (7b)$$

$$\mathcal{I}_{\text{re}}^{\text{ts}} := \left\{ i \in \mathcal{I} \mid |t_i| \geq |s_i| \right\} \quad (7c)$$

where $\epsilon > 0$ is a given, sufficiently small value, introduced as a numerical tolerance in (7a)–(7b). For $\epsilon \rightarrow 0^+$, these sets approach the largest and the smallest possible reduction sets, respectively, namely reducing the most and the least of the variables. Instead, the set defined in (7c) represents an arbitrary trade-off, introduced for the sake of comparison; see Fig. 2.

Remark 3. One could think about performing either the reduction or symmetrization step. However, (i) under Assumption 1, once the system is reduced, the symmetrization step is straightforward, inexpensive and likely effective; (ii) the symmetrization step might be impossible without preliminary reduction, depending on the invertibility of \mathbf{S} .

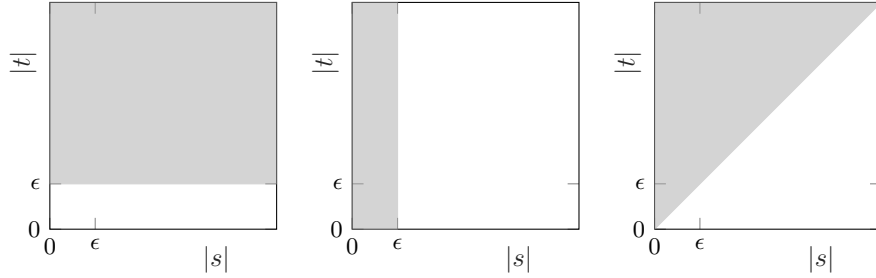


Fig. 2. Reduction sets (7) in the $|s|-|t|$ plane: $\mathcal{I}_{\text{re}}^t$ (left), $\mathcal{I}_{\text{re}}^s$ and $\mathcal{I}_{\text{re}}^{ts}$ (right).

2.1 Reduction

The idea behind the reduction step stems from the observation that problem (1) may be separable, i.e. that it may be possible to compute the value of some variables *a posteriori*, namely once the others are given. In fact, a solution to problem (1) must satisfy

$$\mathbf{Tz} = \mathbf{SCx} + \mathbf{h} \quad (8)$$

where matrix \mathbf{T} is diagonal. Given an index set $\mathcal{I}_{\text{re}} \subseteq \mathcal{I}_{\text{re}}^0$, it is possible to compute z_i , for every $i \in \mathcal{I}_{\text{re}}$, from (8) once the solution vector \mathbf{x} is known. To be sure, let us build matrices $\mathbf{T}_{\text{re}} := \text{diag}(t_i \mid i \in \mathcal{I}_{\text{re}})$ and $\mathbf{T}_{\overline{\text{re}}} := \text{diag}(t_i \mid i \notin \mathcal{I}_{\text{re}})$ and define \mathbf{z}_{re} and $\mathbf{z}_{\overline{\text{re}}}$ the corresponding vectors of unknown variables which can and cannot be reduced, respectively. Then, partitioning the linear system (1) accordingly with these definitions yields:

$$\begin{bmatrix} \mathbf{Q} & \mathbf{A}^\top & \mathbf{C}_{\text{re}}^\top & \mathbf{C}_{\overline{\text{re}}}^\top \\ \mathbf{A} & & & \\ -\mathbf{S}_{\text{re}}\mathbf{C}_{\text{re}} & & \mathbf{T}_{\text{re}} & \\ -\mathbf{S}_{\overline{\text{re}}}\mathbf{C}_{\overline{\text{re}}} & & & \mathbf{T}_{\overline{\text{re}}} \end{bmatrix} \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \\ \mathbf{z}_{\text{re}} \\ \mathbf{z}_{\overline{\text{re}}} \end{pmatrix} = \begin{pmatrix} \mathbf{f} \\ \mathbf{g} \\ \mathbf{h}_{\text{re}} \\ \mathbf{h}_{\overline{\text{re}}} \end{pmatrix} \quad (9)$$

where matrices \mathbf{C}_{re} , $\mathbf{C}_{\overline{\text{re}}}$, \mathbf{S}_{re} , $\mathbf{S}_{\overline{\text{re}}}$ and vectors \mathbf{h}_{re} and $\mathbf{h}_{\overline{\text{re}}}$ are constructed analogously, based on \mathcal{I}_{re} . The matrix \mathbf{T}_{re} is nonsingular, by definition, and then, from (9), one can formally solve for \mathbf{z}_{re} , obtaining

$$\mathbf{z}_{\text{re}} = \mathbf{T}_{\text{re}}^{-1} (\mathbf{S}_{\text{re}}\mathbf{C}_{\text{re}}\mathbf{x} + \mathbf{h}_{\text{re}}) , \quad (10)$$

whose evaluation is straightforward because \mathbf{T}_{re} is diagonal. Plugging (10) back into (9) leads to a smaller linear system, after rearrangements, without reduced variables \mathbf{z}_{re} , namely:

$$\begin{bmatrix} \hat{\mathbf{Q}} & \mathbf{A}^\top & \mathbf{C}_{\overline{\text{re}}}^\top \\ \mathbf{A} & & \\ -\mathbf{S}_{\overline{\text{re}}}\mathbf{C}_{\overline{\text{re}}} & & \mathbf{T}_{\overline{\text{re}}} \end{bmatrix} \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \\ \mathbf{z}_{\overline{\text{re}}} \end{pmatrix} = \begin{pmatrix} \hat{\mathbf{f}} \\ \mathbf{g} \\ \mathbf{h}_{\overline{\text{re}}} \end{pmatrix} \quad (11)$$

where matrix \hat{Q} and vector \hat{f} are defined by:

$$\hat{Q} := Q + C_{\text{re}}^\top T_{\text{re}}^{-1} S_{\text{re}} C_{\text{re}} \quad (12a)$$

$$\hat{f} := f - C_{\text{re}}^\top T_{\text{re}}^{-1} h_{\text{re}} \quad (12b)$$

The larger the set \mathcal{I}_{re} , the more reduced variables, the smaller the obtained linear system (11). In turn, the computation of \hat{Q} may be costly, involving a matrix-matrix multiplication, Eq. 12a. Also, for sparse problems, the fill-up of matrix \hat{Q} may become significant. These drawbacks suggest there might be a trade-off in the reduction step, and hence an optimal reduction strategy, as pointed out in Remark 2.

2.2 Symmetrization

Linear systems with a symmetric matrix can be solved more efficiently, in terms of time and memory. In order to get a symmetric matrix out of (11), it would suffice to left-multiply the rows associated with inequality constraints, namely with z_{re} , by the inverse of $-S_{\text{re}}$. As discussed above, it is $\mathcal{I}_{\text{re}} \subseteq \mathcal{I}_{\text{sy}}^0$, hence the matrix S_{re} is nonsingular, by construction; moreover, its inversion is straightforward, since it is diagonal. Then, the reduced symmetric linear system $\hat{V}\hat{d} = \hat{r}$ reads:

$$\begin{bmatrix} \hat{Q} & A^\top & C_{\text{re}}^\top \\ A & & \\ C_{\text{re}} & -S_{\text{re}}^{-1}T_{\text{re}} & \end{bmatrix} \begin{pmatrix} x \\ y \\ z_{\text{re}} \end{pmatrix} = \begin{pmatrix} \hat{f} \\ g \\ -S_{\text{re}}^{-1}h_{\text{re}} \end{pmatrix} \quad (13)$$

The matrix \hat{V} is symmetric and smaller than V ; the vector \hat{d} collects the unknowns corresponding to optimization variables (x), equality constraints' multipliers (y) and not-reduced inequality constraints' multipliers (z_{re}).

Remark 4. In (12)–(13), the matrix-matrix products $T_{\text{re}}^{-1}S_{\text{re}}$, $S_{\text{re}}^{-1}T_{\text{re}}$ and the matrix-vector products $T_{\text{re}}^{-1}h_{\text{re}}$, $S_{\text{re}}^{-1}h_{\text{re}}$ can be evaluated as entry-wise vector-vector products. In fact, this is possible because matrices S_{re} , S_{re} , T_{re} and T_{re} are diagonal. Furthermore, one can exploit this feature by choosing a specific multiplication ordering, aiming at the lowest possible computational complexity.

3 Numerical results

This Section reports and discusses the results obtained from a MATLAB [17] implementation of Algorithm 1, considering Remark 4. The plain code (as well as a Julia 1.0 and a Python 3.6 implementation) are publicly available [3].

We are interested in comparing the computation time for solving problem (1), through direct methods, with and without the proposed structure exploitation method, see Fig. 1 above. Also, we investigate how it is affected by the problem size $N := n_x + n_a + n_c$ and the relative number of equality and inequality constraints, $\alpha := n_a/n_x$ and $\gamma := n_c/n_x$, respectively.

Algorithm 1 Abstract structure-exploiting linear solver.

Input: $\mathbf{Q}, \mathbf{A}, \mathbf{C}, \mathbf{s}, \mathbf{t}, \mathbf{f}, \mathbf{g}, \mathbf{h}; \epsilon$
Output: $\mathbf{x}, \mathbf{y}, \mathbf{z}$
 $\mathcal{I}_{\text{re}} \leftarrow \mathbf{s}, \mathbf{t}, \epsilon;$ // reduction strategy, Eq. 7
 $\mathbf{C}_{\text{re}}, \mathbf{s}_{\text{re}}, \mathbf{t}_{\text{re}}, \mathbf{h}_{\text{re}}, \mathbf{C}_{\overline{\text{re}}}, \mathbf{s}_{\overline{\text{re}}}, \mathbf{t}_{\overline{\text{re}}}, \mathbf{h}_{\overline{\text{re}}} \leftarrow \mathbf{C}, \mathbf{s}, \mathbf{t}, \mathbf{h}, \mathcal{I}_{\text{re}};$ // partitioning
 $\hat{\mathbf{Q}} \leftarrow \mathbf{Q}, \mathbf{C}_{\text{re}}, \mathbf{t}_{\text{re}}, \mathbf{s}_{\text{re}};$ // Eq. 12a
 $\hat{\mathbf{f}} \leftarrow \mathbf{f}, \mathbf{C}_{\text{re}}, \mathbf{t}_{\text{re}}, \mathbf{h}_{\text{re}};$ // Eq. 12b
 $\hat{\mathbf{V}} \leftarrow \hat{\mathbf{Q}}, \mathbf{A}, \mathbf{C}_{\overline{\text{re}}}, \mathbf{s}_{\overline{\text{re}}}, \mathbf{t}_{\overline{\text{re}}};$ // Eq. 13
 $\hat{\mathbf{r}} \leftarrow \hat{\mathbf{f}}, \mathbf{g}, \mathbf{s}_{\overline{\text{re}}}, \mathbf{h}_{\overline{\text{re}}};$ // Eq. 13
 $\mathbf{x}, \mathbf{y}, \mathbf{z}_{\overline{\text{re}}} \leftarrow \hat{\mathbf{V}}, \hat{\mathbf{r}};$ // linear system
 $\mathbf{z}_{\text{re}} \leftarrow \mathbf{C}_{\text{re}}, \mathbf{s}_{\text{re}}, \mathbf{t}_{\text{re}}, \mathbf{h}_{\text{re}}, \mathbf{x};$ // recovering, Eq. 10
 $\mathbf{z} \leftarrow \mathbf{z}_{\text{re}}, \mathbf{z}_{\overline{\text{re}}}, \mathcal{I}_{\text{re}};$ // assembling

A problem instance consists of matrices $\mathbf{Q}, \mathbf{A}, \mathbf{C}$ and vectors \mathbf{s}, \mathbf{t} (the diagonal of \mathbf{S} and \mathbf{T} , respectively), \mathbf{f}, \mathbf{g} and \mathbf{h} . In the case of full matrices, starting from given values of N, α and γ , an instance is generated as follows:

$$\begin{aligned} n_x &= \left\lceil \frac{N}{1 + \alpha + \gamma} \right\rceil \\ n_a &= \lceil \alpha n_x \rceil \\ n_c &= N - n_x - n_a \\ \bar{Q}_{ij} &\sim \mathcal{N}(0, 1) & i = 1, \dots, n_x, j = 1, \dots, n_x \\ \mathbf{Q} &= \frac{1}{2} (\bar{\mathbf{Q}} + \bar{\mathbf{Q}}^\top) \\ A_{ij} &\sim \mathcal{N}(0, 1) & i = 1, \dots, n_a, j = 1, \dots, n_x \\ C_{ij} &\sim \mathcal{N}(0, 1) & i = 1, \dots, n_c, j = 1, \dots, n_x \\ \rho_i &\sim \sqrt{\mathcal{U}(0, 1)} & i = 1, \dots, n_c \\ \theta_i &\sim \mathcal{U}(0, 2\pi) & i = 1, \dots, n_c \\ s_i &= 1 + \rho_i \cos \theta_i & i = 1, \dots, n_c \\ t_i &= 1 + \rho_i \sin \theta_i & i = 1, \dots, n_c \\ f_i &\sim \mathcal{N}(0, 1) & i = 1, \dots, n_x \\ g_i &\sim \mathcal{N}(0, 1) & i = 1, \dots, n_a \\ h_i &\sim \mathcal{N}(0, 1) & i = 1, \dots, n_c \end{aligned}$$

where $\mathcal{N}(\mu, \sigma)$ denotes the normal continuous probability distribution with mean value μ and standard deviation σ , and $\mathcal{U}(a, b)$ the uniform distribution with support in $[a, b]$. Entries of \mathbf{S} and \mathbf{T} are pairwise coupled in that they are sampled from a disk in the s - t plane, centered in $(1, 1)$ with unitary radius, with uniform probability distribution. This setting is motivated by and mimics the generalized differential of the Fischer-Burmeister function [8]. Both, the direct and the structure-exploiting methods setup the linear system starting from these inputs. Notice that the reduced approach does not build \mathbf{V} nor \mathbf{r} , but their reduced

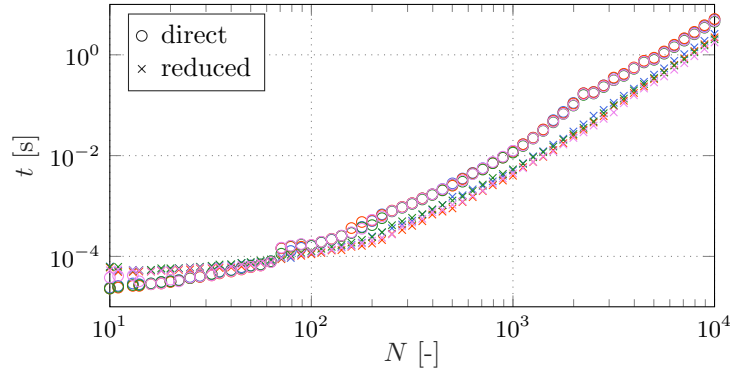


Fig. 3. Execution time: direct approach (t_{dir} , dot) and reduced approach (t_{red} , cross) with reduction set $\mathcal{I}_{\text{re}}^t$; full matrices. Median value.

and symmetric counterpart $\hat{\mathbf{V}}$ and $\hat{\mathbf{r}}$. In our implementation, the direct method builds \mathbf{V} and \mathbf{r} and then adopts the `mldivide` routine to solve $\mathbf{V}\mathbf{d} = \mathbf{r}$. Instead, for the reduced approach (with full matrices), the `linsolve` routine is adopted and explicitly informed that matrix $\hat{\mathbf{V}}$ is symmetric. The problem size N varies between 10 and 10^4 for full and between 10^3 and $2 \cdot 10^4$ for sparse matrices. For each problem size, a set of 100 problem instances are generated (only 10 if $N > 10^4$), checked for ill-conditioning and eventually solved. The composition of constraints is chosen to be $(\alpha, \gamma) \in \{(0, 1), (0, 1.5), (0.5, 1), (0.5, 1.5)\}$ (colored in blue, red, green and violet, respectively). The index sets defined in (7) are adopted and compared, with the tolerance $\epsilon = 10^{-3}$. Sparse matrices are generated in such a way that they approximately have 10 entries for each row; this makes the number of nonzero entries to increase linearly and not quadratically with the problem size N .

The computation time for the direct and reduced case are depicted in Fig. 3, considering full matrices and the (large) index set $\mathcal{I}_{\text{re}}^t$. This gives an idea about the adopted implementation and computing hardware; also, one can guess the computational complexity of the underlying algorithm for solving a linear system. As expected in Section 2, the overhead due to partitioning, reducing and recovering, introduces a break-even point, at around $N = 60$ (for $\mathcal{I}_{\text{re}}^t$ and $\mathcal{I}_{\text{re}}^{\text{ts}}$, but not for $\mathcal{I}_{\text{re}}^s$); hence, the reduced approach is not beneficial only for small-sized problems. This and other considerations can be drawn based on Fig. 4, where it is depicted the (median value of the) ratio of the execution time with the reduced approach, t_{red} , with the different reduction strategies, over the direct one, t_{dir} . Therein, the break-even point corresponds to the unitary ratio; also, the additional computational burden is significant for low values of N . For large N , instead, the ratio decreases to approximately one-half for $\mathcal{I}_{\text{re}}^t$ and $\mathcal{I}_{\text{re}}^{\text{ts}}$, while for $\mathcal{I}_{\text{re}}^s$ it stays around the unit. As one could expect, the reduction set $\mathcal{I}_{\text{re}}^s$ is not as effective as $\mathcal{I}_{\text{re}}^t$ and $\mathcal{I}_{\text{re}}^{\text{ts}}$ because it does not benefit very much from the reduction step, in that it eliminates only few variables. The relative number

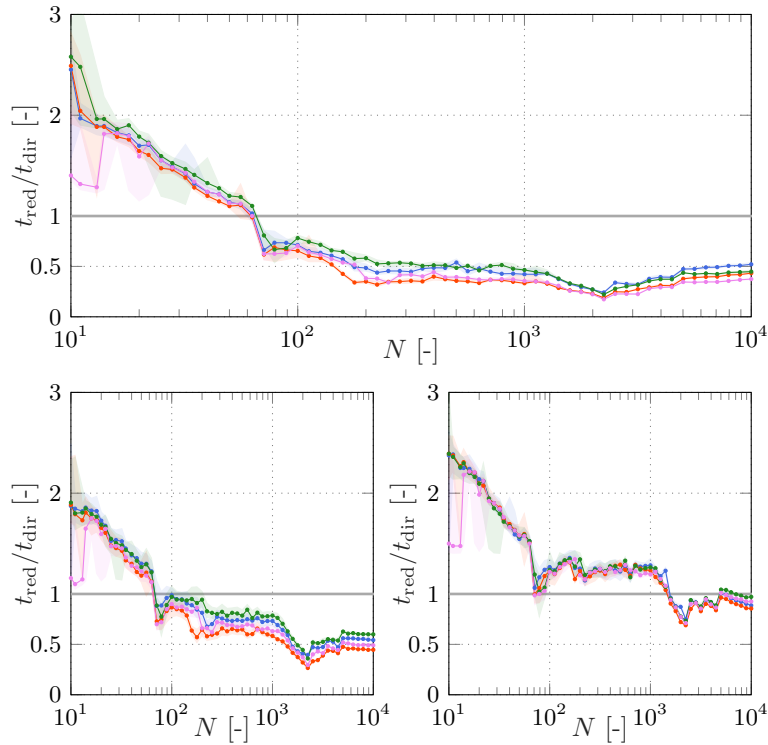


Fig. 4. Ratio of execution time: $\mathcal{I}_{\text{re}}^t$ (top), $\mathcal{I}_{\text{re}}^{\text{ts}}$ (bottom left), $\mathcal{I}_{\text{re}}^s$ (bottom right); full matrices. Median value (line) and 9%-91% quantiles (filled area).

of constraints has an impact on the execution time but it does not drastically affect the overall behaviour (for both, full and sparse case). In fact, all else being equal, either decreasing the number of equalities and increasing the number of inequalities reduces the execution time ratio, meaning that the reduced approach is more effective and worthy for (large) problems with many inequality constraints. For what concerns the case of sparse matrices, similar observations are valid, see Fig. 5. In order to show the distribution of the results obtained from the executed tests, along with the median value, the 9% and 91% quantiles are also reported. For full matrices, Fig. 4, the distribution is relatively narrow, while for sparse matrices, Fig. 5, the results are relatively scattered. Thus, we argue the sparsity pattern greatly affects the computation time. Nevertheless, for relatively large sparse matrices, the ratio $t_{\text{red}}/t_{\text{dir}}$ approaches one-half and promisingly decreases.

These numerical results suggest the set $\mathcal{I}_{\text{re}}^t$ defined in (7a) to be the most effective reduction strategy among those tested. In fact, it generates the smallest linear system and then post-solves the most variables. However, as argued

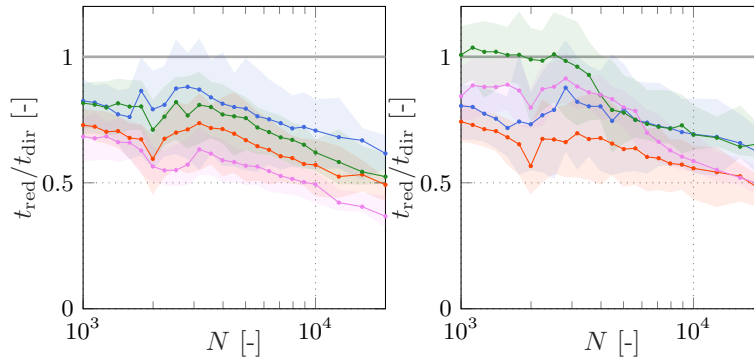


Fig. 5. Ratio of execution time: $\mathcal{I}_{\text{re}}^t$ (top) and $\mathcal{I}_{\text{re}}^{\text{ts}}$ (bottom); sparse matrices. Median value (line) and 9%-91% quantiles (filled area).

in Section 2.1, we claim it might be not the case for (much) larger problem instances, because of the required overhead for the reduction step.

4 Conclusions

This paper proposed and studied a structure-exploiting approach for solving linear systems arising in the context of nonsmooth Newton’s method. The applicability of this method was established under well-posedness of the original problem. Numerical examples showed that the developed approach reduces the computational time, with both, full and sparse linear systems. Some of the tested reduction strategies resulted in halving the execution time.

Analogous ideas apply when an iterative linear solver is of choice, e.g., for very large systems; tailored preconditioners are subject of future research. It remains to assess the effectiveness and the drawbacks of the method when embedded into larger routines for numerical optimization. Moreover, it would be interesting to investigate an optimal reduction strategy, possibly computationally aware.

Acknowledgements

A.D.M. heartily thanks his *marafiki* for the unforgettable memories and wishes them a happy marriage.

References

1. Chen, B., Chen, X., Kanzow, C.: A penalized Fischer-Burmeister NCP-function. *Mathematical Programming* **88**(1), 211–216 (Jun 2000). <https://doi.org/10.1007/PL00011375>
2. Clarke, F.H.: *Optimization and Nonsmooth Analysis*. John Wiley & Sons, New York (1983)

3. De Marchi, A.: Code supporting “Nonsmooth Newton’s method: some structure exploitation” (Nov 2018). <https://doi.org/10.5281/zenodo.1486064>
4. Duff, I.S., Reid, J.K.: MA48 — a Fortran code for direct solution of sparse unsymmetric linear systems of equations. Tech. Rep. RAL-93-072, Rutherford Appleton Laboratory (Oct 1993)
5. Facchinei, F., Fischer, A., Kanzow, C.: On the accurate identification of active constraints. *SIAM Journal on Optimization* **9**(1), 14–32 (1998). <https://doi.org/10.1137/S1052623496305882>
6. Facchinei, F., Kanzow, C.: A nonsmooth inexact Newton method for the solution of large-scale nonlinear complementarity problems. *Mathematical Programming* **76**(3), 493–512 (Mar 1997). <https://doi.org/10.1007/BF02614395>
7. Fischer, A.: A special Newton-type optimization method. *Optimization* **24**, 269–284 (1992). <https://doi.org/10.1080/02331939208843795>
8. Gerdt, M., Kunkel, M.: A nonsmooth Newton’s method for discretized optimal control problems with state and control constraints. *Journal of Industrial & Management Optimization* **4**(2), 247–270 (2008). <https://doi.org/10.3934/jimo.2008.4.247>
9. Gerdt, M., Kunkel, M.: A globally convergent semi-smooth Newton method for control-state constrained DAE optimal control problems. *Computational Optimization and Applications* **48**(3), 601–633 (Apr 2011). <https://doi.org/10.1007/s10589-009-9275-0>
10. Hintermüller, M., Ito, K., Kunisch, K.: The primal-dual active set strategy as a semismooth Newton method. *SIAM Journal on Optimization* **13**(3), 865–888 (2002). <https://doi.org/10.1137/S1052623401383558>
11. Jiang, H.: Global convergence analysis of the generalized Newton and Gauß-Newton methods of the Fischer-Burmeister equation for the complementarity problem. *Mathematics of Operations Research* **24**(3), 529–543 (1999), <http://www.jstor.org/stable/3690647>
12. Laiu, M.P., Tits, A.L.: A constraint-reduced MPC algorithm for convex quadratic programming, with a modified active set identification scheme. *Computational Optimization and Applications* (Mar 2019). <https://doi.org/10.1007/s10589-019-00058-0>
13. Li, X.S.: An overview of SuperLU: Algorithms, implementation, and user interface. *ACM Transactions on Mathematical Software* **31**(3), 302–325 (Sep 2005)
14. Qi, L., Sun, J.: A nonsmooth version of Newton’s method. *Mathematical Programming* **58**(1), 353–367 (Jan 1993). <https://doi.org/10.1007/BF01581275>
15. Schenk, O., Gärtner, K.: Solving unsymmetric sparse systems of linear equations with PARDISO. *Future Generation Computer Systems* **20**(3), 475–487 (2004). <https://doi.org/https://doi.org/10.1016/j.future.2003.07.011>
16. Sun, D., Qi, L.: On NCP-functions. *Computational Optimization and Applications* **13**(1), 201–220 (Apr 1999). <https://doi.org/10.1023/A:1008669226453>
17. The MathWorks, Inc.: MATLAB Release 2017b, Natick, Massachusetts, United States
18. Ulbrich, M.: Nonsmooth Newton-like Methods for Variational Inequalities and Constrained Optimization Problems in Function Spaces. Ph.D. thesis, Technische Universität München (Feb 2002)