



A linear algebra perspective on the random multi-block ADMM: the QP case

Stefano Cipolla¹ · Jacek Gondzio¹

Received: 9 March 2023 / Revised: 21 August 2023 / Accepted: 13 October 2023
© The Author(s) 2023

Abstract

Embedding randomization procedures in the Alternating Direction Method of Multipliers (ADMM) has recently attracted an increasing amount of interest as a remedy to the fact that the direct multi-block generalization of ADMM is not necessarily convergent. Even if, in practice, the introduction of such techniques could *mitigate* the diverging behaviour of the multi-block extension of ADMM, from the theoretical point of view, it can ensure just the *convergence in expectation*, which may not be a good indicator of its robustness and efficiency. In this work, analysing the strongly convex quadratic programming case from a linear algebra perspective, we interpret the block Gauss–Seidel sweep performed by the multi-block ADMM in the context of the inexact Augmented Lagrangian Method. Using the proposed analysis, we are able to outline an alternative technique to those present in the literature which, supported from stronger theoretical guarantees, is able to ensure the convergence of the multi-block generalization of the ADMM method.

Keywords Alternating direction method of multipliers · Inexact augmented Lagrangian method · Randomly shuffled Gauss–Seidel

Mathematics Subject Classification 90C25 · 65K05 · 65F10

1 Introduction

In this work we consider the solution of the Quadratic Programming (QP) problem:

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) := \frac{1}{2} \mathbf{x}^T H \mathbf{x} + \mathbf{g}^T \mathbf{x} \quad (\text{QP})$$

✉ Stefano Cipolla
s.cipolla@soton.ac.uk

✉ Jacek Gondzio
j.gondzio@ed.ac.uk

¹ School of Mathematics, The University of Edinburgh, Peter Guthrie Tait Road, Edinburgh EH93FD, UK

$$\text{s.t. } \mathbf{Ax} = \mathbf{b},$$

where $H \in \mathbb{R}^{d \times d}$ is Symmetric Positive Definite (SPD in short) and $A \in \mathbb{R}^{m \times d}$ ($d \geq m$) has full rank.

Recently, problem (QP) has been widely used as a *sample problem* for the convergence analysis of the n -block generalization of the Alternating Direction Method of Multipliers (ADMM) [6, 11, 22, 49, 55]. In particular, in [11], a counterexample in the form of problem (QP) has been given to show that the direct n -block extension of ADMM is not necessarily convergent when solving non-separable convex minimization problems. This counterexample has motivated a series of very recent works, including [8, 10, 12, 14, 19, 29, 32–35, 41, 42, 44, 53, 54], where the authors analyse modifications of ADMM which ensure its convergence when $n \geq 3$. In particular, in [12, 44, 53] a series of randomization procedures has been introduced which is able to guarantee the convergence in expectation of the n -block generalization of ADMM. Since then such techniques have been proposed as a possible remedy to the fact that the deterministic direct n -block extension of ADMM is not necessarily convergent.

The ADMM [6, 22] was originally proposed in [28] and, in its n -block version, it embeds a n -block Gauss–Seidel (GS) decomposition [5, 27] into each iteration of the Augmented Lagrangian Method (ALM) [36, 50]: the primal variables, partitioned into n blocks, are cyclically updated and then a dual-ascent-type step for the dual variables is performed.

Adopting a purely linear-algebraic approach, in the particular case of problem (QP), ALM and ADMM can be simply interpreted in terms of matrix splitting techniques (see [31, 56]) for the solution of the corresponding Karush–Kuhn–Tucker (KKT) linear system (see Sects. 3 and 6).

Even if in the numerical linear algebra community the study of matrix splitting techniques for the solution of linear systems arising from saddle point problems is a well established line of research (see [2, Sec. 8] for an overview), this connection seems to be only partially exploited in the works [12, 44, 53] and, despite the fact that analogies between ADMM and GS+ALM are apparent, to the best of our knowledge, very few works perform a precise investigation in this direction (even in the simple case when the problem is given by Eq. (QP)).

Indeed, even if it is natural to view ADMM as an approximate version of the ALM, as reported in [22, 23], there were no known results in quantifying this interpretation until the very recent work [15]: here the authors investigate the connection of the block symmetric Gauss–Seidel method [31, Sec. 4.1.1] with the inexact proximal ALM, which represents somehow a different setting from the one investigated here.

Broadly speaking, this work aims to depict a precise picture of the synergies occurring between GS and ALM in order to give rise to ADMM and, in turn, to shed new light on the hidden machinery which controls its convergence.

For the reasons explained above, our starting point is an analysis of the ALM from an inexact point of view and specifically tailored for problem (QP). Indeed, inexact ALMs (iALM) have attracted the attention of many researchers in the last years and we refer to [57, Sec. 1.4] for a very recent literature review. We mention explicitly the works [38, 39, 43, 45], where iALM is analysed for solving linearly constrained

convex programming problems, a very similar framework to the one analysed here. To the best of our knowledge, our approach does not have any evident analogy to the previously mentioned papers.

On the other hand, the connections of the ALM with monotone operators/splitting methods are well understood [21, 51] and, our analysis, resembles this line of research more closely: we use, in essence, a matrix splitting of the augmented KKT matrix of (QP) to represent the ALM/iALM iterations. It is not surprising that, as a result of this line of reasoning, we are able to relate the convergence of ALM/iALM (and their rate of convergence to an ε -accurate primal-dual solution) to the spectral radius ρ of the iteration map of a fixed point problem (see Eq. (9)).

A careful checking of the literature revealed some analogies of our approach with the inexact Uzawa's method [1]. Indeed the ALM method can be interpreted as the Uzawa's method applied to the augmented KKT system of problem (QP) and in the context of the inexact Uzawa's method, it is empirically well documented [25] and theoretically well understood [7, 16–18, 24], that a fixed number of Successive Over-Relaxation (SOR) [26, 58] steps per inner solve (typically 10) is needed in order to reproduce the convergence rate of the exact algorithm.

All the inexactness criteria developed in the previously mentioned works are characterized by a *summability condition* or a *relative error* condition based on the residual previously computed.

A first important by-product of our analysis, is that we are able to prove the convergence of the iALM without imposing any summability condition on the sequence $\{\eta^k\}_k$ which controls the *amount of inexactness* of the iALM at k -th iteration (see Theorem 9) also in the case when the source of inexactness is modelled using a random variable (see Lemma 10). A second important advantage of our approach, is that we are able to give explicit bounds for the rate of convergence of the iALM in relation to the speed characterizing the convergence to zero of the sequence $\{\eta^k\}_k$.

Beyond the previously mentioned advantages of our analysis, we trace the main contribution of this work in the production of an explicit link between the accuracy required to ensure the convergence and the specific solver used to address the minimization step in the ALM, which, in the case of problem (QP), is equivalent to the solution of a SPD linear system. Using explicit error-reduction bounds for the SOR method [47] and its Randomly Shuffled version [48], we are able to prove that the inexactness criterion $\eta^k = R^{k+1}$ ($R < 1$ suitably user-defined), can be satisfied by performing a constant number of iterations (see Theorem 17). Moreover, observing that the GS decomposition is a particular case of the SOR decomposition, we are able to connect the very well known convergence issues [49, 55] of the direct n -block extension of ADMM (and its randomized versions [12, 44, 53]) to the fact that one GS sweep for iALM-step may not be sufficient to ensure enough of the accuracy in the algorithm to deliver convergence. Finally, as an interesting result of our analysis, we are able to propose a *simple* numerical strategy aiming to mitigate, if not to eliminate entirely, the convergence issues of ADMM (see Sect. 6): this proposal, due to its solid theoretical guarantees of convergence, could be considered as a competitive alternative to the techniques introduced to date [12, 44, 53]. We provide some preliminary computational evidence of this fact (see Sect. 7).

1.1 Test problems

In order to showcase the developed theory, in the remainder of this work, we will consider the following test problems (all the numerical results presented are obtained using Matlab® R2020b):

Problem 1 H is the Kernel Matrix associated with the radial basis function for the dataset `heart_scale` from [9] (270 instances, 13 features). In particular, we consider $(H)_{ij} = e^{-\frac{\|x_i - x_j\|}{h^2}}$ with $h = 0.5$ and \mathbf{g} a random vector. For the constraints, we choose $A = \mathbf{e}^T$ where \mathbf{e} is the vector of all ones and $\mathbf{b} = 1$.

Problem 2 Following [11], we consider $H = hI_{3 \times 3}$ with $h = 0.05$ and \mathbf{g} a random vector. For the constraints we consider the matrix

$$A = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 2 \\ 1 & 2 & 2 \end{bmatrix}$$

and \mathbf{b} a random vector ($rank(A) = 3$).

1.2 Notation

In the following, as it is customary in the optimization community, we will use superscripts to denote particular elements of a sequence (scalars, vectors, matrices). In particular, for scalars, this choice could locally clash with the power operation of a scalar but the meaning will be always clear from the context. To denote the whole sequence, instead, we will use subscripts, e.g., $\{\eta^k\}_k \in \mathbb{R}$, $\{\mathbf{x}^k\}_k \in \mathbb{R}^d$ and so on. Given a vector $\mathbf{v} \in \mathbb{R}^d$, $\|\mathbf{v}\|$ denotes the Euclidean norm whereas, given a matrix $H \in \mathbb{R}^{d \times d}$, $\|H\|$ denotes the 2-norm. Moreover, $k_2(H) := \|H\| \|H^{-1}\|$ will be used for the condition number in 2-norm. Finally, in the following, we will freely use a series of standard definitions from linear algebra, e.g., that of minimal polynomial, diagonalizability and similarity, and we refer the interested reader to [37].

2 Augmented Lagrangian and KKT

If we consider the Augmented Lagrangian

$$\mathcal{L}_\beta(\mathbf{x}, \boldsymbol{\mu}) = \frac{1}{2} \mathbf{x}^T H \mathbf{x} + \mathbf{g}^T \mathbf{x} - \boldsymbol{\mu}^T (A \mathbf{x} - \mathbf{b}) + \frac{\beta}{2} \|A \mathbf{x} - \mathbf{b}\|^2,$$

the corresponding KKT conditions are

$$\begin{aligned} \nabla_{\mathbf{x}} \mathcal{L}_\beta(\mathbf{x}, \boldsymbol{\mu}) &= H \mathbf{x} + \mathbf{g} - A^T \boldsymbol{\mu} + \beta A^T A \mathbf{x} - \beta A^T \mathbf{b} = 0 \\ A \mathbf{x} - \mathbf{b} &= 0. \end{aligned}$$

Multiplying by β the second KKT condition, we obtain the system

$$\underbrace{\begin{bmatrix} H_\beta & -A^T \\ \beta A & 0 \end{bmatrix}}_{=: \mathcal{A}} \begin{bmatrix} \mathbf{x} \\ \boldsymbol{\mu} \end{bmatrix} = \underbrace{\begin{bmatrix} \beta A^T \mathbf{b} - \mathbf{g} \\ \beta \mathbf{b} \end{bmatrix}}_{=: \mathbf{q}} \tag{1}$$

where $H_\beta := H + \beta A^T A$. As it will be clear in Sect. 3, the main reason for undergoing the rearrangement of the KKT conditions as in (1), is that we will be able to interpret the Augmented Lagrangian Method of Multipliers (ALM) as a stationary method corresponding to a particular *splitting* of the matrix \mathcal{A} .

Theorem 1 states the existence of a unique solution of problem (1):

Theorem 1 *The matrix \mathcal{A} is invertible for all $\beta > 0$.*

Proof Observe that

$$\mathcal{A} = \begin{bmatrix} H_\beta & 0 \\ \beta A & \beta A H_\beta^{-1} A^T \end{bmatrix} \begin{bmatrix} I - H_\beta^{-1} A^T \\ 0 & I \end{bmatrix}.$$

The non-singularity follows by using the fact that A is of full rank. See also [2, Sec. 3] for different factorizations of saddle point matrices. □

Let us define:

Definition 2 (*ε -accurate primal-dual solution*) We say that $[\mathbf{x}, \boldsymbol{\mu}]^T$ is an ε -accurate primal-dual solution for problem (QP) if

$$\|H\mathbf{x} + \mathbf{g} - A^T \boldsymbol{\mu}\| \leq \varepsilon \text{ and } \|A\mathbf{x} - \mathbf{b}\| \leq \varepsilon.$$

Moreover, if $[\mathbf{x}, \boldsymbol{\mu}]^T$ is a random variable, we say that it is an expected ε -accurate primal-dual solution for problem (QP) if

$$\mathbb{E}(\|H\mathbf{x} + \mathbf{g} - A^T \boldsymbol{\mu}\|) \leq \varepsilon \text{ and } \mathbb{E}(\|A\mathbf{x} - \mathbf{b}\|) \leq \varepsilon.$$

3 The Augmented Lagrangian method of Multipliers (ALM)

The general form of ALM is given by

$$\begin{cases} \mathbf{x}^{k+1} = \min_{\mathbf{x} \in \mathbf{R}^d} \mathcal{L}_\beta(\mathbf{x}, \boldsymbol{\mu}^k) \\ \boldsymbol{\mu}^{k+1} = \boldsymbol{\mu}^k - \beta(A\mathbf{x}^{k+1} - \mathbf{b}), \end{cases}$$

which, for problem (QP), reads as

$$\begin{cases} \mathbf{x}^{k+1} = H_\beta^{-1}(A^T \boldsymbol{\mu}^k + \beta A^T \mathbf{b} - \mathbf{g}) \\ \boldsymbol{\mu}^{k+1} = \boldsymbol{\mu}^k - \beta(A\mathbf{x}^{k+1} - \mathbf{b}) \end{cases}. \tag{2}$$

It is important to observe that the iterates $[\mathbf{x}^{k+1}, \boldsymbol{\mu}^{k+1}]^T$ produced by (2) are dual feasible, i.e.,

$$0 = \nabla_{\mathbf{x}} \mathcal{L}_{\beta}(\mathbf{x}^{k+1}, \boldsymbol{\mu}^k) = \nabla_{\mathbf{x}} f(\mathbf{x}^{k+1}) - A^T \boldsymbol{\mu}^{k+1} = H\mathbf{x}^{k+1} + \mathbf{g} - A^T \boldsymbol{\mu}^{k+1}.$$

It is well known that ALM can be derived applying the Proximal Point Method to the dual of problem (QP), see [51, Sec. 6.1], but in this particular case can be also recast in an operator splitting framework (see [51, Sec. 7], [21]): indeed, the ALM scheme can be interpreted as a fixed point iteration obtained from a splitting decomposition for the KKT linear algebraic system (1) (see [30, 56] and [2, Sec. 8]). Writing

$$\mathcal{A} = \begin{bmatrix} H_{\beta} & 0 \\ \beta A & I \end{bmatrix} - \begin{bmatrix} 0 & A^T \\ 0 & I \end{bmatrix},$$

we can write Eq. (2) as

$$\begin{bmatrix} \mathbf{x}^{k+1} \\ \boldsymbol{\mu}^{k+1} \end{bmatrix} = \underbrace{\begin{bmatrix} H_{\beta}^{-1} & 0 \\ -\beta A H_{\beta}^{-1} & I \end{bmatrix} \begin{bmatrix} 0 & A^T \\ 0 & I \end{bmatrix}}_{=: G_{\beta}} \begin{bmatrix} \mathbf{x}^k \\ \boldsymbol{\mu}^k \end{bmatrix} + \underbrace{\begin{bmatrix} H_{\beta}^{-1} & 0 \\ -\beta A H_{\beta}^{-1} & I \end{bmatrix}}_{=: F_{\beta}} \underbrace{\begin{bmatrix} \beta A^T \mathbf{b} - \mathbf{g} \\ \beta \mathbf{b} \end{bmatrix}}_{\mathbf{q}},$$

i.e., as a fixed point iteration of the form

$$\begin{bmatrix} \mathbf{x}^{k+1} \\ \boldsymbol{\mu}^{k+1} \end{bmatrix} = G_{\beta} \begin{bmatrix} \mathbf{x}^k \\ \boldsymbol{\mu}^k \end{bmatrix} + F_{\beta} \mathbf{q}.$$

The following Theorem 3 (see [13, Sec. 2] for a similar result) is the cornerstone to prove the convergence of the ALM (see Eq. (2)) and its inexact version (see Eq. (8)).

Theorem 3 *The eigenvalues of G_{β} are s.t. $\lambda \in [0, 1)$ for all $\beta > 0$ and, moreover, $\rho(G_{\beta}) \rightarrow 0$ for $\beta \rightarrow \infty$.*

Proof Let us observe that $(\lambda, [\mathbf{u}, \mathbf{v}]^T)$ is an eigenpair of G_{β} if and only if

$$\begin{aligned} A^T \mathbf{v} &= \lambda H_{\beta} \mathbf{u} \\ (1 - \lambda) \mathbf{v} &= \lambda \beta A \mathbf{u}. \end{aligned} \tag{3}$$

The proof is structured into three parts.

Part 1: If λ is an eigenvalue of G_{β} , then $\lambda \neq 1$.

By contradiction suppose that $\lambda = 1$, then from (3) we have the condition

$$\begin{bmatrix} H_{\beta} & -A^T \\ \beta A & 0 \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} = 0,$$

which leads to an absurd since \mathcal{A} is invertible for $\beta > 0$ (see Theorem 1).

Part 2: If $(\lambda, [\mathbf{u}, \mathbf{v}]^T)$ is an eigenpair of G_{β} , then $\mathbf{u} \neq 0$.

By contradiction, if $\mathbf{u} = 0$, then from the second equation in (3), we obtain $(1-\lambda)\mathbf{v} = 0$ and hence an absurd using Part 1.

Part 3: If $\mathbf{v} = 0$, multiplying by \mathbf{u}^T the first equation in (3), we obtain $\lambda\mathbf{u}^T H_\beta \mathbf{u} = 0$, which leads to $\lambda = 0$ since H_β is SPD.

If $\mathbf{v} \neq 0$, from (3), we obtain

$$\lambda(1-\lambda)\frac{\mathbf{u}^T H_\beta \mathbf{u}}{\mathbf{u}^T \mathbf{u}} = \lambda\beta\frac{\mathbf{u}^T A^T A \mathbf{u}}{\mathbf{u}^T \mathbf{u}}. \tag{4}$$

If in Eq. (4) $\frac{\mathbf{u}^T A^T A \mathbf{u}}{\mathbf{u}^T \mathbf{u}} = 0$, reasoning as before and using Part 1, we obtain $\lambda = 0$. Instead, if in Eq. (4) we have $\frac{\mathbf{u}^T A^T A \mathbf{u}}{\mathbf{u}^T \mathbf{u}} \neq 0$, we obtain $\lambda = 0$ or $\lambda = \frac{\mathbf{u}^T H \mathbf{u}}{\mathbf{u}^T H_\beta \mathbf{u}} < 1$, which completes the proof observing that, in this case, $\lambda = \frac{\mathbf{u}^T H \mathbf{u}}{\mathbf{u}^T H_\beta \mathbf{u}} \rightarrow 0$ if $\beta \rightarrow \infty$ since $\mathbf{u}^T A^T A \mathbf{u} \neq 0$. □

Lemma 4 *The matrix G_β is diagonalizable.*

Proof Let us start from observing that

$$G_\beta = \begin{bmatrix} 0 & H_\beta^{-1} A^T \\ 0 & I - \beta A H_\beta^{-1} A^T \end{bmatrix}. \tag{5}$$

The proof is divided into two parts.

Part 1: We prove that the matrix $I - \beta A H_\beta^{-1} A^T$ is invertible.

To prove this fact, it is enough to prove that $\beta A H_\beta^{-1} A^T$ does not have unitary eigenvalues. Using Woodbury formula and defining $C := (I + \beta A H^{-1} A^T)^{-1}$, we have

$$\begin{aligned} \beta A H_\beta^{-1} A^T &= \beta A H^{-1} A^T (I - C \beta A H^{-1} A^T) \Rightarrow \\ \beta A H_\beta^{-1} A^T \mathbf{x} = \mathbf{x} &\Leftrightarrow \beta A H^{-1} A^T C \mathbf{x} = \mathbf{x}. \end{aligned}$$

This follows by observing that $\beta A H^{-1} A^T C$ and $C^{\frac{1}{2}} \beta A H^{-1} A^T C^{\frac{1}{2}}$ are similar and that

$$\lambda(C^{\frac{1}{2}} \beta A H^{-1} A^T C^{\frac{1}{2}}) \subset (0, 1).$$

Part 2: We prove that the minimal polynomial of G_β factorizes in distinct linear factors.

The proof of this fact follows by observing that the minimal polynomials of the blocks on the diagonal of G_β , namely the null matrix and the symmetric matrix $I - \beta A H_\beta^{-1} A^T$, factorize in distinct linear factors since they are diagonalizable (see [37, Cor 3.3.10]). Moreover, since the matrix $I - \beta A H_\beta^{-1} A^T$ is invertible, they do not have common factors. Hence, the product of such minimal polynomials (which coincides with the lowest common multiple (*lcm*)) is the minimal polynomial of the whole matrix.

Indeed, this last implication holds for generic block upper triangular matrices. To prove that, let us consider a block upper triangular matrix F with b diagonal blocks F_{ii} , $i = 1, \dots, b$. Let us denote, moreover, by $m_i(x)$ the minimal polynomials of the blocks and $m(x)$ the minimal polynomial of the whole matrix F . We have $lcm(m_i(x)) | m(x)$ because $m(F_{ii}) = 0$. Moreover, by direct computation, one can check that, defining $s(x) := \prod_{i=1}^n m_i(x)$, it holds $s(G) = 0$. If the polynomials $m_i(x)$ are pairwise relatively prime (i.e., they do not have common factors), then $s(x) = lcm(m_i(x))$ and hence $s(x) = m(x)$ since $lcm(m_i(x)) | m(x)$.

The final proof of statement, i.e., the diagonalizability of G_β , follows by observing that, if the minimal polynomial of a given matrix factorizes in distinct linear factors, then the matrix is diagonalizable (see, once more, [37, Cor 3.3.10]). \square

Remark 1 (Non unitary step length) The framework presented until now allows to consider also the case of non-unitary dual steps. Indeed, considering the splitting

$$\mathcal{A} = \begin{bmatrix} H_\beta & 0 \\ \beta A & \frac{1}{\gamma} I \end{bmatrix} - \begin{bmatrix} 0 & A^T \\ 0 & \frac{1}{\gamma} I \end{bmatrix},$$

it is easy to see that the corresponding ALM-type update is

$$\begin{cases} \mathbf{x}^{k+1} = \mathbf{x}^{k+1} = H_\beta^{-1}(A^T \boldsymbol{\mu}^k + \beta A^T \mathbf{b} - \mathbf{g}) \\ \boldsymbol{\mu}^{k+1} = \boldsymbol{\mu}^k - \gamma \beta (A \mathbf{x}^{k+1} - \mathbf{b}). \end{cases}$$

Moreover, using the techniques from Theorem 3 and Lemma 4, it can be proved that its convergence and rate of convergence depend on the spectral radius of

$$I - \gamma \beta A H_\beta^{-1} A^T.$$

Hence, the choice of the parameter γ could be used, in principle, to further improve the rate of convergence of ALM. In the following we consider $\gamma = 1$.

Lemma 5 *There exists a constant $M \equiv M(G_\beta) \geq 1$ s.t. $\|G_\beta^k\| \leq M \rho(G_\beta)^k$ for all $k \in \mathbb{N}$.*

Proof Using Lemma 4, since G_β is diagonalizable, we have $G_\beta^k = X \Lambda^k X^{-1}$, and hence

$$\|G_\beta^k\| \leq \underbrace{\|X\| \|X^{-1}\|}_{=: M} \|\Lambda^k\| \leq M \rho(G_\beta)^k. \tag{6}$$

\square

Definition 6 In the following, $[\bar{\mathbf{x}}, \bar{\boldsymbol{\mu}}]^T$ denotes the unique solution of linear system (1) (see Theorem 1 for existence and uniqueness). Moreover, we define, $\rho_\beta := \rho(G_\beta) := \max_\lambda \{|\lambda(G_\beta)|\}$, $\mathbf{e}^k := \begin{bmatrix} \mathbf{x}^k - \bar{\mathbf{x}} \\ \boldsymbol{\mu}^k - \bar{\boldsymbol{\mu}} \end{bmatrix}$, $\mathbf{d}^k := \mathcal{A} \begin{bmatrix} \mathbf{x}^k \\ \boldsymbol{\mu}^k \end{bmatrix} - \mathbf{q}$.

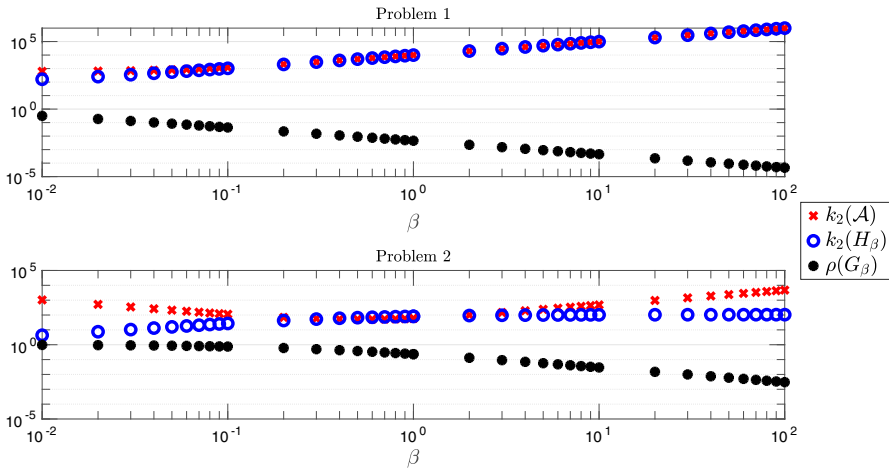


Fig. 1 Behaviour of $k_2(\mathcal{A})$, $k_2(H_\beta)$, $\rho(G_\beta)$ for different values of β (logarithmic scale on both axes). Using Lemma 7 the ALM method (2) converges faster for smaller values of $\rho(G_\beta)$.

Lemma 7 *The ALM in (2) converges for all $\beta > 0$. Moreover, we have for all $k \in \mathbb{N}$,*

$$\|\mathbf{e}^k\| \leq \|\mathbf{e}^0\| M \rho_\beta^k$$

and

$$\|\mathbf{d}^k\| \leq \|\mathcal{A}\| \|\mathcal{A}^{-1}\| \|\mathbf{d}^0\| M \rho_\beta^k.$$

Proof From direct computation, we have

$$\begin{aligned} \mathbf{e}^k &= G_\beta^k \mathbf{e}^0, \\ \mathbf{d}^k &= \mathcal{A} G_\beta^k \mathcal{A}^{-1} \mathbf{d}^0, \end{aligned}$$

where we used $\mathcal{A} \mathbf{e}^k = \mathbf{d}^k$. Thesis follows by passing to the norms and using Lemma 5. □

In Fig. 1 we report the behaviour of the condition number in 2-norm of the matrices \mathcal{A} , H_β (respectively $k_2(\mathcal{A})$, $k_2(H_\beta)$) and the spectral radius ρ_β for different values of β . The results obtained in Fig. 1 confirm the statement regarding ρ_β in Theorem 3, i.e., ρ_β decreases when β increases. And indeed, using Lemma 7, we can observe that the convergence of ALM can be consistently sped-up by increasing the value of β , which corresponds to a decrease of ρ_β . On the other hand, the eventual speed-up resulting from considering *large* values for β comes at the cost of solving an increasingly ill-conditioned linear system involving H_β (see the first equation in (2) and the behaviour of $k_2(H_\beta)$ in Fig. 1). Indeed, when β is large, the matrix H_β is dominated by the term $\beta A^T A$ (see [2, Sec. 8.1] and references therein for more details) and, if $A^T A$

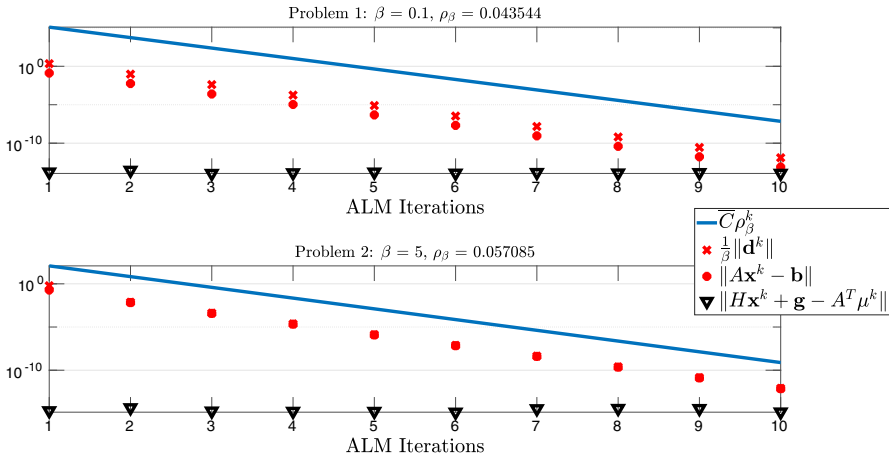


Fig. 2 Behaviour of the quantities analysed in Lemma 8 (logarithmic scale on y-axis)

is singular, the condition number of the matrix H_β progressively degrades when β increases (see the behaviour of $k_2(H_\beta)$ for Problem 1 in the upper panel of Fig. 1).

The following Lemma 8 states the worst case complexity of ALM.

Lemma 8 *The ALM in (2) requires $O(\log_{\rho_\beta} \varepsilon)$ iterations to produce an ε -accurate primal-dual solution.*

Proof Observe that we have

$$\|Ax^k - b\| \leq \frac{1}{\beta} \|d^k\| \leq \frac{1}{\beta} \|A\| \|A^{-1}\| \|d^0\| M \rho_\beta^k,$$

where in the last inequality we used Lemma 7. Since, as observed at the beginning of this section, the iterates $[x^k, \mu^k]^T$ produced by the ALM are dual feasible, we have $\|Hx^k + g - A^T \mu^k\| \equiv 0$. Hence, defining $\bar{C} := \frac{1}{\beta} \|A\| \|A^{-1}\| \|d^0\| M$, we obtain that $k \geq \log_{\rho_\beta} (\varepsilon / \bar{C})$ iterations of the ALM are sufficient to deliver an ε -accurate primal-dual solution. \square

In Fig. 2, we show the behaviour of the quantities involved in the proof of Lemma 8 (the notation used in the legend is consistent with that used in Lemma 8 except for the fact that the numerical value of the constant \bar{C} used there is normalized using M , i.e., in Fig. 2 we report $\bar{C} \equiv \bar{C}/M$). As Lemma 8 states and Fig. 8 shows, the function $\bar{C} \rho_\beta^k$ is an upper bound for the quantity $\|Ax^k - b\|$. In this example, in order to further highlight the dependence of ρ_β on β , we choose different values of β ($\beta = 0.1$ and $\beta = 5$) such that, for Problem 1 and Problem 2, we obtain $\rho_\beta \approx 0.05$. Let us point out that the results reported in Fig. 2 are obtained solving the linear system in (2) using a high accuracy (a direct method using Matlab’s “backslash” operator) and, since the iterates must be dual feasible, the residuals $\|Hx^k + g - A^T \mu^k\|$ are close to the machine precision.

4 Inexact ALM (iALM)

In this section we study in detail the iALM for problem (QP). The reader may see [57, Sec. 1.4] for a recent survey on this subject. In particular, we assume that the first equation in (2) is not solved exactly, i.e., \mathbf{x}^{k+1} is such that

$$H_\beta \mathbf{x}^{k+1} - (A^T \boldsymbol{\mu}^k + \beta A^T \mathbf{b} - \mathbf{g}) = \mathbf{r}^k. \tag{7}$$

In our framework, the iALM read as

$$\begin{cases} \mathbf{x}^{k+1} = H_\beta^{-1} (A^T \boldsymbol{\mu}^k + \beta A^T \mathbf{b} - \mathbf{g} + \mathbf{r}^k), \\ \boldsymbol{\mu}^{k+1} = \boldsymbol{\mu}^k - \beta (A \mathbf{x}^{k+1} - \mathbf{b}) \end{cases}, \tag{8}$$

and (8) can be alternatively written as the following inexact fixed point iteration (see [4] and [46, Sec 12.2] for more details on this topic):

$$\begin{bmatrix} \mathbf{x}^{k+1} \\ \boldsymbol{\mu}^{k+1} \end{bmatrix} = \underbrace{\begin{bmatrix} H_\beta^{-1} & 0 \\ -\beta A H_\beta^{-1} & I \end{bmatrix} \begin{bmatrix} 0 & A^T \\ 0 & I \end{bmatrix}}_{=:F_\beta} \begin{bmatrix} \mathbf{x}^k \\ \boldsymbol{\mu}^k \end{bmatrix} + \underbrace{\begin{bmatrix} H_\beta^{-1} & 0 \\ -\beta A H_\beta^{-1} & I \end{bmatrix}}_{=:F_\beta} \underbrace{\begin{bmatrix} \beta A^T \mathbf{b} - \mathbf{g} + \mathbf{r}^k \\ \beta \mathbf{b} \end{bmatrix}}_{=: \mathbf{q}^k}. \tag{9}$$

On the contrary of what was observed for the exact ALM (see the beginning of Sect. 3), the iterates produced by (8) are not dual feasible since

$$0 \neq \mathbf{r}^k = \nabla_{\mathbf{x}} \mathcal{L}_\beta(\mathbf{x}^{k+1}, \boldsymbol{\mu}^k) = H \mathbf{x}^{k+1} + \mathbf{g} - A^T \boldsymbol{\mu}^{k+1},$$

i.e., the error \mathbf{r}^k introduced in the solution of the first equation in (8) can be interpreted as a measure of the violation of the dual feasibility condition.

In Sect. 5.1.2 we will consider the point \mathbf{x}^{k+1} in (7) as a result of a *randomized* procedure and, for this reason, we are going to present this section assuming that $\{\mathbf{r}^k\}_k$ in (9) is a sequence of random variables (and hence all the generated $\{[\mathbf{x}^k, \boldsymbol{\mu}^k]^T\}_k$ are random variables). Moreover, all the results presented here can be easily restated in a *deterministic* framework substituting the “almost sure (a.s.) convergence” with “convergence” and not considering the “expectation operator”. For a review of the probabilistic concepts we use in the following see [52, Ch. 2].

The following Theorem 9 addresses the convergence of the iALM using the inexact fixed point formulation in (9) under the condition that $\|\mathbf{r}^k\|$, i.e., the error introduced in the solution of the first equation in (8), converges a.s. to zero.

Theorem 9 *Let $\beta > 0$. If $\lim_{j \rightarrow \infty} \|\mathbf{r}^j\| = 0$ a.s., then the iALM in (8) converges a.s. to the solution of the linear system (1) and the following inequalities hold a.s. for every $k \in \mathbb{N}$:*

$$\|\mathbf{e}^k\| \leq M \rho_\beta^k \|\mathbf{e}^0\| + M \|F_\beta\| \sum_{j=0}^{k-1} \rho_\beta^{k-1-j} \|\mathbf{r}^j\|,$$

$$\|\mathbf{d}^k\| \leq \|\mathcal{A}\|\|\mathcal{A}\|^{-1}\|\mathbf{d}^0\|M\rho_\beta^k + M\|\mathcal{A}\|\|F_\beta\| \sum_{j=0}^{k-1} \rho_\beta^{k-1-j} \|\mathbf{r}^j\|. \tag{10}$$

Proof If $[\bar{\mathbf{x}}, \bar{\boldsymbol{\mu}}]^T$ is a solution of (1), then it satisfies the fixed point equation

$$\begin{bmatrix} \bar{\mathbf{x}} \\ \bar{\boldsymbol{\mu}} \end{bmatrix} = \begin{bmatrix} H_\beta^{-1} & 0 \\ -\beta AH_\beta^{-1} & I \end{bmatrix} \begin{bmatrix} 0 & A^T \\ 0 & I \end{bmatrix} \begin{bmatrix} \bar{\mathbf{x}} \\ \bar{\boldsymbol{\mu}} \end{bmatrix} + \begin{bmatrix} H_\beta^{-1} & 0 \\ -\beta AH_\beta^{-1} & I \end{bmatrix} \begin{bmatrix} \beta A^T \mathbf{b} - \mathbf{g} \\ \beta \mathbf{b} \end{bmatrix}. \tag{11}$$

Subtracting (11) from (9) we obtain

$$\mathbf{e}^k = G_\beta \mathbf{e}^{k-1} + F_\beta \begin{bmatrix} \mathbf{r}^{k-1} \\ 0 \end{bmatrix} \text{ a.s.}$$

and hence

$$\mathbf{e}^k = G_\beta^k \mathbf{e}^0 + \sum_{j=0}^{k-1} G_\beta^{k-1-j} F_\beta \begin{bmatrix} \mathbf{r}^j \\ 0 \end{bmatrix} \text{ a.s.} \tag{12}$$

Passing to the norms in (12) and using Lemma 5, we have

$$\|\mathbf{e}^k\| \leq \rho_\beta^k M \|\mathbf{e}^0\| + M \|F_\beta\| \sum_{j=0}^{k-1} \rho_\beta^{k-1-j} \|\mathbf{r}^j\| \text{ a.s.} \tag{13}$$

The a.s. convergence to zero of $\{\|\mathbf{e}^k\|\}_k$ follows from (13) observing that, if $\lim_{k \rightarrow \infty} \|\mathbf{r}^k\| = 0$ a.s., then

$$\lim_{k \rightarrow \infty} \sum_{j=0}^{k-1} \rho_\beta^{k-1-j} \|\mathbf{r}^j\| = 0 \text{ a.s.}$$

(this is a particular case of the Toeplitz Lemma, see [46, Exercise 12.2-3] for the deterministic case, [40] and references therein for the probabilistic case). The second part of the statement follows by observing that

$$\|\mathbf{d}^k\| = \|\mathcal{A}\mathbf{e}^k\| \leq \|\mathcal{A}\|\|\mathbf{e}^k\| \text{ a.s.}$$

and that $\|\mathbf{e}^0\| \leq \|\mathcal{A}\|^{-1}\|\mathbf{d}^0\|$. □

Lemma 10 *Suppose $\mathbb{E}(\|\mathbf{r}^j\|) \leq R^{j+1}$ for all $j \in \mathbb{N}$ and $R < 1$. Then the iterates of the iALM in (8) converge a.s. to the solution of the linear system (1). Moreover, if $R < \rho_\beta$, then $O(\log_{\rho_\beta} \varepsilon)$ iterations are sufficient to produce an expected ε -accurate primal-dual solution; else, if $\rho_\beta \leq R$, then $O(\log_R \varepsilon)$ iterations are sufficient (given that ε is sufficiently small).*

Proof If $\mathbb{E}(\|\mathbf{r}^j\|) \leq R^{j+1}$ for all $j \in \mathbb{N}$, then $\sum_{j=0}^{\infty} \mathbb{E}(\|\mathbf{r}^j\|) < \infty$ and hence, using [52, Th. 2.1.3], we have $\lim_{j \rightarrow \infty} \|\mathbf{r}^j\| = 0$ a.s. Using now Theorem 9, we have that $\|\mathbf{d}^k\|$ converges a.s. to zero.

Using Eq. (10) and the hypothesis $\mathbb{E}(\|\mathbf{r}^j\|) \leq R^{j+1}$, we have

$$\mathbb{E}(\|\mathbf{d}^k\|) \leq \|\mathcal{A}\| \|\mathcal{A}\|^{-1} \|\mathbf{d}^0\| M \rho_\beta^k + M \|\mathcal{A}\| \|F_\beta\| \sum_{j=0}^{k-1} \rho_\beta^{k-1-j} R^{j+1}. \tag{14}$$

Let us observe, moreover, that

$$\begin{aligned} &\mathbb{E}(\|H\mathbf{x}^k + \mathbf{g} - A^T \boldsymbol{\mu}^k\| - \|\beta A^T(A\mathbf{x}^k - \mathbf{b})\|) \\ &\leq \mathbb{E}(\|H_\beta \mathbf{x}^k - A^T \boldsymbol{\mu}^k + \mathbf{g} - \beta A^T \mathbf{b}\|) \leq \mathbb{E}(\|\mathbf{d}^k\|), \end{aligned}$$

and hence

$$\mathbb{E}(\|H\mathbf{x}^k + \mathbf{g} - A^T \boldsymbol{\mu}^k\|) \leq \mathbb{E}(\|\mathbf{d}^k\|) + \|A^T\| \mathbb{E}(\|\mathbf{d}^k\|) \leq C_1 \mathbb{E}(\|\mathbf{d}^k\|), \tag{15}$$

where we defined $C_1 := (1 + \|A^T\|)$ and used the fact that $\|\beta(A\mathbf{x}^k - \mathbf{b})\| \leq \|\mathbf{d}^k\|$ a.s.

Case $R < \rho_\beta$. Using (14), we have

$$\mathbb{E}(\|\mathbf{d}^k\|) \leq \|\mathcal{A}\| \|\mathcal{A}\|^{-1} \|\mathbf{d}^0\| M \rho_\beta^k + \rho_\beta^k M \|\mathcal{A}\| \|F_\beta\| \frac{R}{\rho_\beta} \sum_{j=0}^{k-1} \left(\frac{R}{\rho_\beta}\right)^j \leq C_2 \rho_\beta^k,$$

where $C_2 := \max\{M \|\mathcal{A}\| \|\mathcal{A}\|^{-1} \|\mathbf{d}^0\|, M \frac{\frac{R}{\rho_\beta} \|\mathcal{A}\| \|F_\beta\|}{1 - \frac{R}{\rho_\beta}}\}$.

Moreover, using the above inequality, we have also

$$\mathbb{E}(\|A\mathbf{x}^k - \mathbf{b}\|) \leq \frac{1}{\beta} \mathbb{E}(\|\mathbf{d}^k\|) \leq \frac{1}{\beta} C_2 \rho_\beta^k,$$

and hence, using (15) and defining $\bar{C} := \max\{C_1 C_2, \frac{1}{\beta} C_2\}$, we obtain that $k \geq \log_{\rho_\beta}(\varepsilon/\bar{C})$ iterations of iALM are sufficient to produce an expected ε -accurate primal-dual solution.

Case $\rho_\beta \leq R$. Using (14), we have

$$\mathbb{E}(\|\mathbf{d}^k\|) \leq \|\mathcal{A}\| \|\mathcal{A}\|^{-1} \|\mathbf{d}^0\| M R^k + R^k k M \|\mathcal{A}\| \|F_\beta\| \leq C_2 R^k k,$$

where $C_2 := \max\{M \|\mathcal{A}\| \|\mathcal{A}\|^{-1} \|\mathbf{d}^0\|, M \|\mathcal{A}\| \|F_\beta\|\}$. Let us observe that, in this case, we have

$$\mathbb{E}(\|A\mathbf{x}^k - \mathbf{b}\|) \leq \frac{1}{\beta} \mathbb{E}(\|\mathbf{d}^k\|) \leq \frac{1}{\beta} C_2 R^k k,$$

and hence, using (15) and defining $\bar{C} := \max\{C_1 C_2, \frac{1}{\beta} C_2\}$, we obtain that to produce an expected ε -accurate primal-dual solution it suffices to perform $k + \log_R k \geq \log_R(\varepsilon/\bar{C})$ iterations of iALM. The last part of the statement follows by observing that $\lim_{k \rightarrow \infty} \frac{k + \log_R k}{k} = 1$. \square

Before concluding this section, let us state the following Corollary 11, which will be used later:

Corollary 11 *Suppose $\mathbb{E}(\|\mathbf{r}^j\|) \leq R^{j+1}$ for all $j \in \mathbb{N}$ and $R < 1$. If $R > \rho_\beta$, then there exists a constant $L \equiv L(M, \mathcal{A}, \rho_\beta, R)$ s.t.*

$$\frac{\|\beta(A\mathbf{x}^k - \mathbf{b})\|}{R^k} \leq L < \infty \quad \text{a.s. for every } k \in \mathbb{N}, \tag{16}$$

and hence, we have

$$\mathbb{E}\left(\frac{\|\beta(A\mathbf{x}^k - \mathbf{b})\|}{R^k}\right) \leq L \quad \text{for every } k \in \mathbb{N}. \tag{17}$$

Proof Using (10) we have

$$\begin{aligned} \frac{\|\beta(A\mathbf{x}^k - \mathbf{b})\|}{R^k} &\leq \frac{\|\mathbf{d}^k\|}{R^k} \\ &\leq M\|\mathcal{A}\|\|\mathcal{A}^{-1}\|\|\mathbf{d}^0\| \left(\frac{\rho_\beta}{R}\right)^k + M\|\mathcal{A}\|\|F_\beta\| \sum_{j=0}^{k-1} \left(\frac{\rho_\beta}{R}\right)^{k-1-j}, \end{aligned}$$

from which thesis follows by observing that $\sum_{j=0}^{k-1} \left(\frac{\rho_\beta}{R}\right)^{k-1-j} \leq \frac{1}{1-\frac{\rho_\beta}{R}}$ for all k . \square

5 The solution of the linear system

In this section, given $[\mathbf{x}^k, \boldsymbol{\mu}^k]$, we suppose that the linear system

$$H_\beta \mathbf{x} = (A^T \boldsymbol{\mu}^k + \beta A^T \mathbf{b} - \mathbf{g}), \tag{18}$$

is solved using an iterative solver. Despite the fact that any iterative solver can be used for the solution of the SPD system in (18), we will focus our attention only on the block Successive Over-Relaxation method (SOR) [26, 58] or its Randomly Shuffled version (RSSOR) [48]. Indeed, these choices will allow us to clearly interpret the Random n -block ADMM as an iALM, see Sect. 6. Since \mathbf{r}^k in the first equation of (8) is the (possibly deterministic) residual associated to the linear system (18), i.e.,

$$H_\beta \mathbf{x}^{k+1} - (A^T \boldsymbol{\mu}^k + \beta A^T \mathbf{b} - \mathbf{g}) = \mathbf{r}^k,$$

one would be tempted to think that the *increasing accuracy condition* for the a.s. convergence to zero of the expected residual \mathbf{r}^k in Lemma 10, i.e., $\mathbb{E}(\|\mathbf{r}^k\|) \leq R^{k+1}$,

requires that the expected number of iterations of the chosen iterative solver increases when the iterates of iALM proceed. In this section we will show that this is not the case if $R > \rho_\beta$. For the remaining of this work let us define

$$\chi^k := A^T \mu^k + \beta A^T \mathbf{b} - \mathbf{g},$$

and $\{\eta^k\}_k \rightarrow 0$ as the *forcing sequence* such that $\mathbb{E}(\|\mathbf{r}^k\|) \leq \eta^k$ for all $k \in \mathbb{N}$. We use, moreover, the following inequalities: given $B \in \mathbb{R}^{d \times d}$ SPD, if we order the eigenvalues of B as $\lambda_1(B) \geq \dots \geq \lambda_d(B)$, then

$$\lambda_d(B) \|\mathbf{x}\|_B^2 \leq \|B\mathbf{x}\|^2 \leq \lambda_1(B) \|\mathbf{x}\|_B^2 \quad \text{for all } \mathbf{x} \in \mathbb{R}^d \tag{19}$$

and

$$\lambda_d(B) \|\mathbf{x}\|^2 \leq \|B^{1/2}\mathbf{x}\|^2 \leq \lambda_1(B) \|\mathbf{x}\|^2 \quad \text{for all } \mathbf{x} \in \mathbb{R}^d. \tag{20}$$

For the sake of completeness, before presenting our results, we deliver a brief survey on the block SOR method which is based on [31, 48, 56].

5.1 A brief survey on SOR and randomly shuffled SOR

Let $B \in \mathbb{C}^{d \times d}$. Consider the linear system

$$B\mathbf{y} = \chi. \tag{21}$$

We can express the matrix B as the sum of block-matrices $B = D - L - U$ where

$$D := \begin{bmatrix} B_{1,1} & & & \\ & B_{2,2} & & \\ & & \ddots & \\ & & & B_{n,n} \end{bmatrix}, \quad L := - \begin{bmatrix} 0_{1,1} & 0 & 0 & 0 \\ B_{2,1} & 0_{2,2} & 0 & \vdots \\ \vdots & \ddots & \ddots & 0 \\ B_{n,1} & B_{n,2} & \dots & 0_{n,n} \end{bmatrix},$$

$$U := - \begin{bmatrix} 0_{1,1} & B_{1,2} & & B_{1,n} \\ 0 & 0_{2,2} & \ddots & \vdots \\ \vdots & & \ddots & B_{n-1,n} \\ 0 & 0 & \dots & 0_{n,n} \end{bmatrix}. \tag{22}$$

Let us suppose now that the block-diagonal matrix D is invertible. The fixed point problem corresponding to Eq. (21) can be written as

$$(D - \omega L)\mathbf{y} = ((1 - \omega)D + \omega U)\mathbf{y} + \omega \chi$$

and the SOR method is defined as

$$\mathbf{y}^{j+1} = (D - \omega L)^{-1}((1 - \omega)D + \omega U)\mathbf{y}^j + \omega(D - \omega L)^{-1}\chi. \tag{23}$$

The Gauss–Seidel (GS) method is recovered for $\omega = 1$.

It is important to point out at this stage that, to interpret the block Gauss–Seidel sweep performed by the multi-block ADMM in the context of the inexact Augmented Lagrangian Method as in Sect. 6, we will need the results contained in this section only in the particular case of $\omega = 1$ (which corresponds precisely to the Gauss–Seidel case). On the other hand, we prefer to state all the theory in the more general framework of SOR ($0 < \omega < 2$) as the results presented in Sect. 6 hold also for relaxation parameters different from $\omega = 1$. Despite the fact that the detailed study of such generalizations of the multi-block ADMM falls out of the scope of this work, it is important to note that they might be of great practical interest due to the enhanced rate of convergence of SOR w.r.t. Gauss–Seidel when suitably selected relaxation parameters are chosen.

Observe that Eq. (23) can be written alternatively as

$$\mathbf{y}^{j+1} = (I - \omega D^{-1}L)^{-1}((1 - \omega)I + \omega D^{-1}U)\mathbf{y}^j + \omega(I - \omega D^{-1}L)^{-1}D^{-1}\boldsymbol{\chi}, \tag{24}$$

and for this reason, usually, the *point successive over-relaxation matrix* is defined as

$$\mathcal{L}_\omega := (I - \omega D^{-1}L)^{-1}((1 - \omega)I + \omega D^{-1}U).$$

The following Corollary of the Ostrowski–Reich Theorem states the convergence of the block SOR iteration:

Corollary 12 [56, Cor. 3.14] *Let $B \in \mathbb{C}^{n \times n}$ and D, L, U be defined as in (22). If D is positive definite, then the block SOR method in (23) is convergent for all \mathbf{y}^0 if and only if $0 < \omega < 2$ and B is positive definite.*

In this work we are going to deal just with symmetric matrices and, for this reason, we denote the factor U in (22) with L^T . It is worth noting, moreover, that using the equality $(1 - \omega)D + \omega L^T = (D - \omega L) - \omega B$, we can further rewrite the SOR iteration in (23) as

$$\mathbf{y}^{j+1} = (I - \omega(D - \omega L)^{-1}B)\mathbf{y}^j + \omega(D - \omega L)^{-1}\boldsymbol{\chi}. \tag{25}$$

In [48], a Randomly Shuffled version of SOR (RSSOR) has been introduced and studied: it is obtained considering P^j as a random permutation matrix (with uniform distribution and independent from the current guess \mathbf{y}^j) and applying the SOR splitting to the linear system $P^j B P^{jT} P^j \mathbf{x} = P^j \boldsymbol{\chi}$, i.e., considering

$$P^j B P^{jT} = D_{P^j} - L_{P^j} - L_{P^j}^T.$$

The RSSOR is defined as

$$\mathbf{y}^{j+1} = (I - \omega P^{jT}(D_{P^j} - \omega L_{P^j})^{-1}P^j B)\mathbf{y}^j + \omega P^{jT}(D_{P^j} - \omega L_{P^j})^{-1}P^j \boldsymbol{\chi}. \tag{26}$$

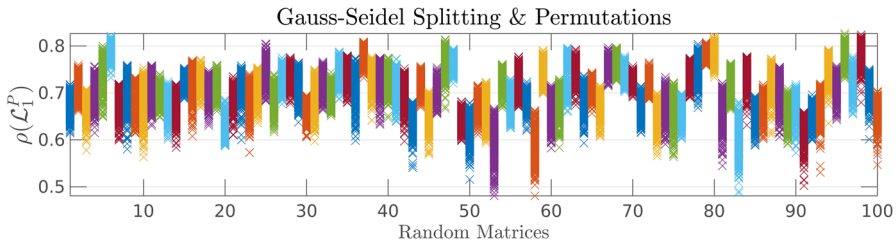


Fig. 3 $\rho(\mathcal{L}_1^P)$ for 100 random matrices of the form $B = R^T R + I \in \mathbb{R}^{7 \times 7}$

Moreover, let us observe that after defining $Q_{Pj} := \omega P^j T (D_{Pj} - \omega L_{Pj})^{-1} P^j$, (26) can be written as a function of the random variables P^1, \dots, P^j , i.e.,

$$y^{j+1} = \prod_{\ell=0}^j (I - Q_{P^\ell} B) y^0 + \sum_{i=0}^j \left(\prod_{\ell=i+1}^j (I - Q_{P^\ell}) \right) Q_{P^i} \chi, \tag{27}$$

where we set $(\prod_{\ell=i+1}^j (I - Q_{P^\ell})) := I$ if $\ell + 1 > j$.

Before concluding this section, let us point out that the main idea connected with RSSOR is related to the fact that, although the spectral distribution of the matrix PBP^T does not depend on any particular permutation matrix P , the spectrum of the lower triangular part $D_P - L_P$ does depend on it. As a result, also the spectral radius of the matrix $\mathcal{L}_\omega^P := (I - \omega P^T (D_P - \omega L_P)^{-1} P B)$ is affected by the particular choice of P . To further highlight the aforementioned dependence and to strengthen the intuition of the reader in this regard, in Fig. 3 we report $\rho(\mathcal{L}_1^P)$ for all the permutation matrices P and for 100 randomly generated matrices of the form $B = R^T R + I \in \mathbb{R}^{7 \times 7}$ (R is generated using the Matlab’s function “rand”).

5.1.1 Rate of convergence of SOR

This section is based on [48]. If B is SPD and is partitioned as in (22), the linear system in (21) can be transformed as

$$D^{-1/2} B D^{-1/2} D^{1/2} y = D^{-1/2} \chi \tag{28}$$

(D is SPD since B is SPD) and hence the coefficient matrix can be decomposed as

$$D^{-1/2} B D^{-1/2} = I - D^{-1/2} L D^{-1/2} - (D^{-1/2} L D^{-1/2})^T, \tag{29}$$

where $D^{-1/2} L D^{-1/2}$ and $(D^{-1/2} L D^{-1/2})^T$ are, respectively, strictly lower triangular and strictly upper triangular. For the above explained reasons, in this section we will suppose that $B = I - L - L^T$.

Observe, moreover, that the SOR method applied to the system in (28) with the splitting (29) coincides exactly with (24) and hence, the fact that in this section we suppose that the diagonal of B is the identity, is expected to simplify the presentation.

The following Theorem 13 gives a precise bound for the rate of convergence of the SOR method:

Theorem 13 [48, Th. 1] *Let B a SPD matrix, then the SOR method (25) converges for $0 < \omega < 2$ in the energy norm associated with B according to*

$$\|\bar{\mathbf{y}} - \mathbf{y}^j\|_B^2 \leq \left(1 - \frac{(2 - \omega)\omega\lambda_1(B)}{\left(1 + \frac{1}{2} \lfloor \log_2(2d) \rfloor \omega\lambda_1(B)\right)^2 k_2(B)}\right)^j \|\bar{\mathbf{y}} - \mathbf{y}^0\|_B^2. \quad (30)$$

The rate of convergence stated in (30) depends on the dimension of the problem d and this feature is not desirable for large scale problems.

One of the main advantages of the RSSOR consists in the fact that the expected error reduction factor is independent from the dimension of the problem, as stated in the following:

Theorem 14 [48, Th. 4] *The expected squared energy norm error of the RSSOR iteration converges exponentially with the bound*

$$\mathbb{E}(\|\bar{\mathbf{y}} - \mathbf{y}^j\|_B^2) \leq \left(1 - \frac{(2 - \omega)\omega\lambda_1(B)}{(1 + \omega\lambda_1(B))^2 k_2(B)}\right)^j \|\bar{\mathbf{y}} - \mathbf{y}^0\|_B^2. \quad (31)$$

for any $\omega \in (0, 2)$.

As already pointed out, Eq. (31) does not exhibit any dependence on the dimension of the problem and, for this reason, the Randomly Shuffled versions of SOR should be considered for large scale problems. Moreover, the following corollary addresses the convergence of the iterates to the solution of the linear system:

Corollary 15 $\lim_{j \rightarrow \infty} \|\bar{\mathbf{y}} - \mathbf{y}^j\|_B^2 = 0$ a.s.

Proof Using (31), we have that $\sum_{j=0}^{\infty} \mathbb{E}(\|\bar{\mathbf{y}} - \mathbf{y}^j\|_B^2) < \infty$. This follows by applying [52, Th. 2.1.3]. \square

5.1.2 Using SOR in iALM

We are ready to analyse the behaviour of SOR method in the framework of the iALM (8). In particular, we are going to present our results for the RSSOR method (see Eq. (26)), but analogous techniques/results apply/hold for the non-randomized version (25). This choice is mainly driven by the reasons of timeliness: in the next Sect. 6 we are able to interpret the recently introduced Randomized ADMM (RADMM) as a particular case of iALM where the linear system (18) is solved (inexactly) using RSSOR with $\omega = 1$ (which will be denoted, in the following, as Randomly Shuffled Gauss–Seidel (RSGS)). For this reason, in this section, we apply the results presented in Sect. 4 in the *probabilistic* form considering $\{\mathbf{r}^k\}_k$ and $\{\mathbf{x}^k, \boldsymbol{\mu}^k\}^T_k$ as sequences of random variables.

Of course, the same results as presented here hold, with *simple* modifications, for the deterministic ADMM and the classical GS method.

In order to use the rate of convergence stated in (31), we write $H_\beta = D - L - L^T$ and transform the linear system in (18) as follows:

$$D^{-1/2} H_\beta D^{-1/2} D^{1/2} \mathbf{x} = D^{-1/2} \boldsymbol{\chi}^k. \tag{32}$$

Let us define $\tilde{H}_\beta = D^{-1/2} H_\beta D^{-1/2}$, $\tilde{\boldsymbol{\chi}}^k := D^{-1/2} \boldsymbol{\chi}^k$, $\tilde{\mathbf{x}} := D^{1/2} \mathbf{x}$.

Consider, moreover, the random variable

$$\mathbb{E} \left(\|\tilde{\mathbf{x}}^{k+1} - \tilde{\mathbf{x}}^{k+1,j}\|_{\tilde{H}_\beta}^2 \mid \begin{bmatrix} \mathbf{x}^k \\ \boldsymbol{\mu}^k \end{bmatrix} \right),$$

where $\tilde{H}_\beta \tilde{\mathbf{x}}^{k+1} = \tilde{\boldsymbol{\chi}}^k$ and $\{\tilde{\mathbf{x}}^{k+1,j}\}_j$ is the random sequence generated by RSSOR method in (26) to approximate $\tilde{\mathbf{x}}^{k+1}$, i.e., the solution of problem (32).

The following Lemma 16 will be useful to state the main result of this section:

Lemma 16 *Let us suppose that the RSSOR in Eq. (26) is used for the solution of the linear system (32) with $\mathbf{y}^0 = D^{1/2} \mathbf{x}^k =: \tilde{\mathbf{x}}^{k+1,0}$. If the random variable $(P^{k+1,0}, \dots, P^{k+1,j})$ is independent from $\{[\mathbf{x}^k, \boldsymbol{\mu}^k]^T\}_k$ for every $j, k \in \mathbb{N}$ (beyond the standard assumptions required on the $P^{k+1,j}$ by RSSOR), then*

$$\mathbb{E}(\|\tilde{\mathbf{x}}^{k+1} - \tilde{\mathbf{x}}^{k+1,j}\|_{\tilde{H}_\beta}) \leq \left(1 - \frac{(2 - \omega)\omega\lambda_1(\tilde{H}_\beta)}{(1 + \omega\lambda_1(\tilde{H}_\beta))^2 k_2(\tilde{H}_\beta)} \right)^{j/2} \mathbb{E}(\|\tilde{\mathbf{x}}^{k+1} - \tilde{\mathbf{x}}^{k+1,0}\|_{\tilde{H}_\beta}). \tag{33}$$

Proof Let us observe that, using (27), we can write

$$\|\tilde{\mathbf{x}}^{k+1} - \tilde{\mathbf{x}}^{k+1,j}\|_{\tilde{H}_\beta}^2 = g \left((P^{k+1,0}, \dots, P^{k+1,j-1}), \begin{bmatrix} \mathbf{x}^k \\ \boldsymbol{\mu}^k \end{bmatrix} \right),$$

where g is a deterministic function.

Using the fact that, if the random variable Y is independent from X (see *Freezing Lemma*, [20, Example 5.1.5]), it holds

$$\mathbb{E}(g(Y, X) \mid X) = \mathbb{E}(g(Y, x))_{|x=X},$$

and using (31), we have

$$\begin{aligned} & \mathbb{E} \left(\|\tilde{\mathbf{x}}^{k+1} - \tilde{\mathbf{x}}^{k+1,j}\|_{\tilde{H}_\beta}^2 \mid \begin{bmatrix} \mathbf{x}^k \\ \boldsymbol{\mu}^k \end{bmatrix} \right) \\ & \leq \left(1 - \frac{(2 - \omega)\omega\lambda_1(\tilde{H}_\beta)}{(1 + \omega\lambda_1(\tilde{H}_\beta))^2 k_2(\tilde{H}_\beta)} \right)^j \|\tilde{\mathbf{x}}^{k+1} - \tilde{\mathbf{x}}^{k+1,0}\|_{\tilde{H}_\beta}^2 \text{ a.s.} \end{aligned}$$

Moreover, using the conditional Jensen’s Inequality in the left hand-side of the previous equation (see [3, Th. 34.4]) and then passing the square root, we have

$$\begin{aligned} & \mathbb{E} \left(\|\tilde{\mathbf{x}}^{k+1} - \tilde{\mathbf{x}}^{k+1,j}\|_{\tilde{H}_\beta} \left\| \begin{bmatrix} \mathbf{x}^k \\ \boldsymbol{\mu}^k \end{bmatrix} \right\| \right) \\ & \leq \left(1 - \frac{(2 - \omega)\omega\lambda_1(\tilde{H}_\beta)}{(1 + \omega\lambda_1(\tilde{H}_\beta))^2 k_2(\tilde{H}_\beta)} \right)^{j/2} \|\tilde{\mathbf{x}}^{k+1} - \tilde{\mathbf{x}}^{k+1,0}\|_{\tilde{H}_\beta} \text{ a.s.} \end{aligned}$$

This follows by considering the expectation on both sides of the above inequality and using the properties of the conditional expectation [3, Th. 34.4]. \square

We are now ready to state the following Theorem 17 which summarizes the properties of the iALM in (8) when each sub-problem is solved using RSSOR:

Theorem 17 *Let $\{\eta^k\}_k = R^{k+1}$ with $R > \rho_\beta$. Define*

$$\bar{j}^{(k)} := \min\{j : \mathbb{E}(\|\mathbf{r}^{k,j}\|) \leq \eta^k\}, \tag{34}$$

where $\{\mathbf{r}^{k,j} := H_\beta \mathbf{x}^{k+1,j} - \boldsymbol{\chi}^k\}_j$ is the sequence of random residuals generated by RSSOR initialized using $\mathbf{x}^{k+1,0} = \mathbf{x}^k$. Then, there exists $\bar{j} \in \mathbb{N}$ such that $\bar{j} \geq \bar{j}^{(k)}$ for all k .

Moreover, an expected ε -accurate primal-dual solution of problem (QP) can be obtained in $O(\log_R \varepsilon)$ iALM iterations.

Proof Using (19) in (33) and since the expectation is a linear function, we have

$$\begin{aligned} & \mathbb{E}(\|\tilde{H}_\beta \tilde{\mathbf{x}}^{k+1,j} - \tilde{\boldsymbol{\chi}}^k\|) \\ & \leq \left(1 - \frac{(2 - \omega)\omega\lambda_1(\tilde{H}_\beta)}{(1 + \omega\lambda_1(\tilde{H}_\beta))^2 k_2(\tilde{H}_\beta)} \right)^{j/2} \sqrt{k_2(\tilde{H}_\beta)} \mathbb{E}(\|\tilde{H}_\beta \tilde{\mathbf{x}}^{k+1,0} - \tilde{\boldsymbol{\chi}}^k\|) \end{aligned}$$

and hence, using (20),

$$\begin{aligned} & \mathbb{E}(\|H_\beta \mathbf{x}^{k+1,j} - \boldsymbol{\chi}^k\|) \\ & \leq \left(1 - \frac{(2 - \omega)\omega\lambda_1(\tilde{H}_\beta)}{(1 + \omega\lambda_1(\tilde{H}_\beta))^2 k_2(\tilde{H}_\beta)} \right)^j \sqrt{k_2(\tilde{H}_\beta)k_2(D^{-1})} \mathbb{E}(\|H_\beta \mathbf{x}^k - \boldsymbol{\chi}^k\|), \end{aligned}$$

where we defined $\mathbf{x}^{k+1,j} := D^{-1/2} \tilde{\mathbf{x}}^{k+1,j}$ for $j \geq 1$. If in the above equation we use the definition of $\mathbf{r}^{k+1,j}$, we have

$$\begin{aligned} & \mathbb{E}(\|\mathbf{r}^{k+1,j}\|) \\ & \leq \left(1 - \frac{(2 - \omega)\omega\lambda_1(\tilde{H}_\beta)}{(1 + \omega\lambda_1(\tilde{H}_\beta))^2 k_2(\tilde{H}_\beta)} \right)^{j/2} \sqrt{k_2(\tilde{H}_\beta)k_2(D^{-1})} \mathbb{E}(\|H_\beta \mathbf{x}^k - \boldsymbol{\chi}^k\|^2), \end{aligned}$$

and hence, defining

$$j^{(k)} := \left\lceil \log \left(1 - \frac{(2-\omega)\omega\lambda_1(\tilde{H}_\beta)}{(1+\omega\lambda_1(\tilde{H}_\beta))^2 k_2(\tilde{H}_\beta)} \right) \frac{2\eta^k}{\sqrt{k_2(\tilde{H}_\beta)k_2(D^{-1})\mathbb{E}(\|H_\beta \mathbf{x}^k - \boldsymbol{\chi}^k\|)}} \right\rceil,$$

it holds $\mathbb{E}(\|\mathbf{r}^{k,j^{(k)}}\|) \leq \eta^k$. Observe, moreover, that using the second equation in (8), we have

$$\mathbb{E}(\|H_\beta \mathbf{x}^k - \boldsymbol{\chi}^k\|) = \mathbb{E}(\|\mathbf{r}^{k-1} + \beta A^T(A\mathbf{x}^k - \mathbf{b})\|) \leq \mathbb{E}(\|\mathbf{r}^{k-1}\|) + \|A^T\| \mathbb{E}(\|\beta(A\mathbf{x}^k - \mathbf{b})\|),$$

and hence using the hypothesis $\eta^k = R^{k+1}$ and Eq. (17), we are able to state the existence of a constant $C > 0$ such that

$$\frac{2R^{k+1}}{\sqrt{k_2(\tilde{H}_\beta)k_2(D^{-1})\mathbb{E}(\|H_\beta \mathbf{x}^k - \boldsymbol{\chi}^k\|)}} \geq C \quad \text{for all } k.$$

We obtain

$$\bar{j} := \left\lceil \log \left(1 - \frac{(2-\omega)\omega\lambda_1(\tilde{H}_\beta)}{(1+\omega\lambda_1(\tilde{H}_\beta))^2 k_2(\tilde{H}_\beta)} \right) C \right\rceil \geq j^{(k)} \quad \text{for all } k. \tag{35}$$

From (35), we obtain the first part of the statement observing that $j^{(k)} \geq \bar{j}^{(k)}$ for all k . The last part of the statement follows by observing that with this choice of η^k the hypotheses of Lemma 10 are satisfied. \square

In the upper panels of Fig. 4 we report the quantities analysed in the proof of Lemma 10 (also in this case the notation used in the legend is consistent with that used in Lemma 10 except for the fact that the numerical constant \bar{C} used there is normalized using M , i.e., in Fig. 4 we report $\bar{C} \equiv \bar{C}/M$). The expectations $\mathbb{E}(\|A\mathbf{x}^k - \mathbf{b}\|)$, $\mathbb{E}(\|H\mathbf{x}^k + \mathbf{g} - A^T\boldsymbol{\mu}^k\|)$ and $\mathbb{E}(\|\mathbf{d}^k\|)$ are approximated using the empirical mean over 15 iALM simulations, whereas, for each fixed k and j , $\mathbb{E}(\|\mathbf{r}^{k,j}\|)$ is approximated using the empirical mean of $\mathbb{E}(\mathbb{E}(\|\mathbf{r}^{k,j}\| \mid \begin{bmatrix} \mathbf{x}^k \\ \boldsymbol{\mu}^k \end{bmatrix}))$ over 15 trajectories for $[\mathbf{x}^k, \boldsymbol{\mu}^k]^T$ and 15 simulations of the RSGS step. In the lower panels, we report, for each iALM step and for each simulation, the box-plots of the obtained $\bar{j}^{(k)}$ (see Eq. (34)). As Theorem 17 states and Fig. 4 confirms, $\bar{j}^{(k)}$ shows a *bounded-from-above* behaviour for all the iALM iterations (the choice of the parameters β and R is reported on top of the figure).

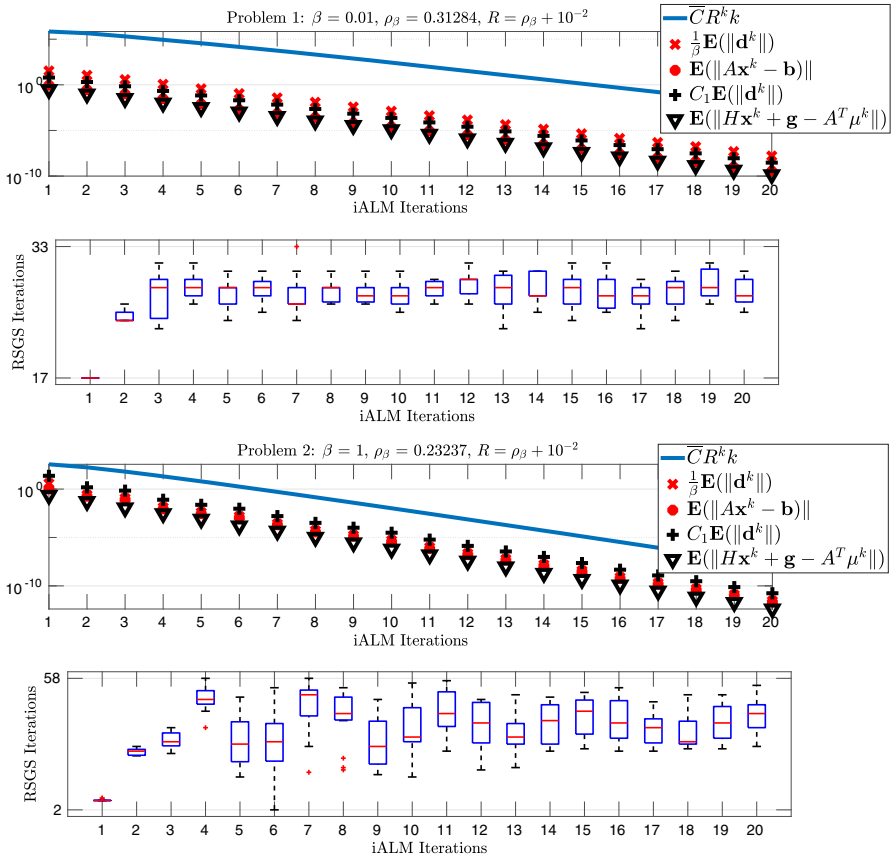


Fig. 4 Upper panels: Behaviour of the quantities analysed in Lemma 10 (logarithmic scale on y-axis) approximated using the empirical mean over 15 simulations of iALM. Lower panels: Box-plots of the $\bar{J}^{(k)}$'s (see Eq. (34)) obtained in each simulation of iALM when RSSOR is used for the solution of (18) using $\{\eta^k\}_k$ and $\{x^{k+1,0}\}_k$ as in Theorem 17

6 Interpreting (random)ADMM as an iALM

Given a block partition of x , i.e., $x = [x_{d_1}, \dots, x_{d_n}]^T$ with $d_1 + \dots + d_n = d$, the n -block ADMM (see [11] and references therein) is defined as

$$\begin{cases} x_{d_1}^{k+1} := \arg \min_{x_{d_1} \in \mathbb{R}^{d_1}} \mathcal{L}_\beta([x_{d_1}, x_{d_2}^k, \dots, x_{d_n}^k]^T, \mu^k), \\ \vdots \\ x_{d_n}^{k+1} := \arg \min_{x_{d_n} \in \mathbb{R}^{d_n}} \mathcal{L}_\beta([x_{d_1}^{k+1}, x_{d_2}^{k+1}, \dots, x_{d_n}]^T, \mu^k), \\ \mu^{k+1} := \mu^k - \beta(Ax^{k+1} - b). \end{cases} \tag{36}$$

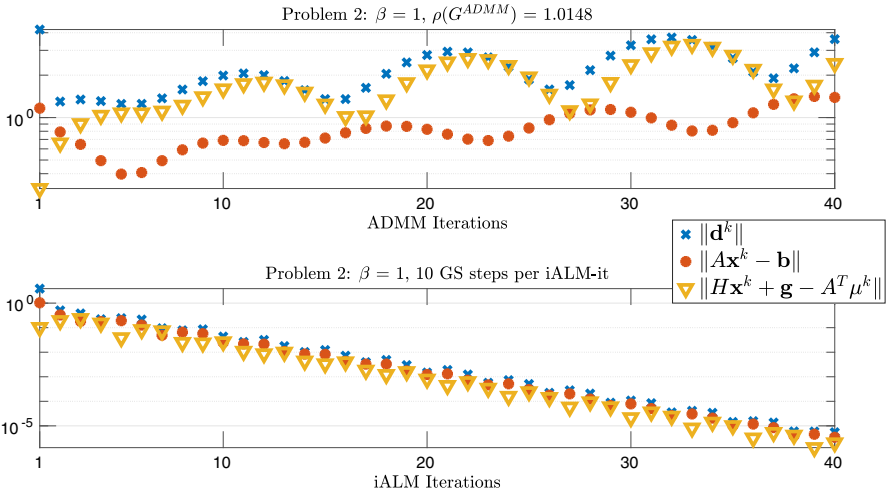


Fig. 5 ADMM vs iALM&GS for Problem 2 (logarithmic scale on y-axis)

If we apply the iterative method in (36) to solve problem (QP), splitting H_β as $H_\beta = D - L - L^T$, it is possible to re-write (36) in compact form (see [12, 53]):

$$\begin{bmatrix} \mathbf{x}^{k+1} \\ \boldsymbol{\mu}^{k+1} \end{bmatrix} = \underbrace{\begin{bmatrix} D - L & 0 \\ \beta A & I \end{bmatrix}^{-1} \begin{bmatrix} L^T & A^T \\ 0 & I \end{bmatrix}}_{=:G^{ADMM}} \begin{bmatrix} \mathbf{x}^k \\ \boldsymbol{\mu}^k \end{bmatrix} + \begin{bmatrix} D - L & 0 \\ \beta A & I \end{bmatrix}^{-1} \begin{bmatrix} \beta A^T \mathbf{b} - \mathbf{g} \\ \beta \mathbf{b} \end{bmatrix}. \quad (37)$$

Since Eq. (37) can be written alternatively as

$$\begin{cases} \mathbf{x}^{k+1} = (D - L)^{-1} L^T \mathbf{x}^k + (D - L)^{-1} (A^T \boldsymbol{\mu}^k + \beta A^T \mathbf{b} - \mathbf{g}) \\ \boldsymbol{\mu}^{k+1} := \boldsymbol{\mu}^k - \beta (A \mathbf{x}^{k+1} - \mathbf{b}), \end{cases} \quad (38)$$

we can observe that the first equation in (38) is precisely one step of the SOR method with $\omega = 1$ (see Eq. (23)), i.e., ADMM performs exactly one GS iteration for the solution of the linear system $H_\beta \mathbf{x} = A^T \boldsymbol{\mu}^k + \beta A^T \mathbf{b} - \mathbf{g}$. Let us point out that in [11] it has been proved that the n -block extension of ADMM is not always convergent since there exist examples where the spectral radius of G^{ADMM} in Eq. (37) satisfies $\rho(G^{ADMM}) > 1$. The analysis performed in Sects. 4 and 5 reveals a simple strategy to remedy this: performing *more* steps of the GS iteration to fulfil the requirements needed on the residuals will ensure convergence. Indeed, as proved in Sect. 5 (deterministic case), a constant number of iterations of SOR per iALM-step is sufficient to guarantee that the produced residuals satisfy the sufficient conditions for convergence. To further underpin the previous claim, in Fig. 5, we report the behaviour of $\|\mathbf{d}^k\|$, $\|A\mathbf{x}^k - \mathbf{b}\|$ and $\|H\mathbf{x}^k + \mathbf{g} - A^T \boldsymbol{\mu}^k\|$ for ADMM and for iALM&GS where, at each inner iteration, 10 GS sweeps are performed. For the particular case of Problem 2 when $\beta = 1$ and all the blocks have size one, we have $\rho(G^{ADMM}) = 1.0148 > 1$ and the ADMM is not convergent (see the upper panel in Fig. 5). On the contrary, performing more than

one GS sweep (lower panel of Fig. 5) is enough to observe a convergent behaviour of all residuals.

Exactly the same observation can be made for the RADMM [12, 53]: this method is obtained considering a block permutation matrix P^k which selects the order for solving the block-equations and then splitting the matrix $P^k H_\beta P^{kT}$ as

$$P^k H_\beta P^{kT} = D_{P^k} - L_{P^k} - L_{P^k}^T \tag{39}$$

(the random permutation matrix is selected independently from the iterate \mathbf{x}^k and uniformly at random among all possible block-permutation matrices). In more details, if we consider the iterative method

$$\left\{ \begin{array}{l} \text{select a permutation } \sigma \text{ of } \{1, \dots, n\} \text{ uniformly at random independently from } \mathbf{x}^k, \\ \mathbf{x}_{d_{\sigma(1)}}^{k+1} := \arg \min_{\mathbf{x}_{d_{\sigma(1)}} \in \mathbf{R}^{d_{\sigma(1)}}} \mathcal{L}_\beta([\mathbf{x}_{d_{\sigma(1)}}^k, \mathbf{x}_{d_{\sigma(2)}}^k, \dots, \mathbf{x}_{d_{\sigma(n)}}^k]^T, \boldsymbol{\mu}^k), \\ \vdots \\ \mathbf{x}_{d_{\sigma(n)}}^{k+1} := \arg \min_{\mathbf{x}_{d_{\sigma(n)}} \in \mathbf{R}^{d_{\sigma(n)}}} \mathcal{L}_\beta([\mathbf{x}_{d_{\sigma(1)}}^{k+1}, \mathbf{x}_{d_{\sigma(2)}}^{k+1}, \dots, \mathbf{x}_{d_{\sigma(n)}}^k]^T, \boldsymbol{\mu}^k), \\ \boldsymbol{\mu}^{k+1} := \boldsymbol{\mu}^k - \beta(A\mathbf{x}^{k+1} - \mathbf{b}) \end{array} \right. \tag{40}$$

to solve problem (QP), using the splitting (39), we can write (40) in the fixed point form

$$\begin{aligned} \begin{bmatrix} \mathbf{x}^{k+1} \\ \boldsymbol{\mu}^{k+1} \end{bmatrix} &= \underbrace{\begin{bmatrix} P_k^T & 0 \\ 0 & I \end{bmatrix} \begin{bmatrix} D_{P_k} - L_{P_k} & 0 \\ \beta A P_k^T & I \end{bmatrix}^{-1} \begin{bmatrix} L_{P_k}^T P_k & P_k A^T \\ 0 & I \end{bmatrix}}_{=: G_\beta^{P_k}} \begin{bmatrix} \mathbf{x}^k \\ \boldsymbol{\mu}^k \end{bmatrix} \\ &+ \begin{bmatrix} P_k^T & 0 \\ 0 & I \end{bmatrix} \begin{bmatrix} D_{P_k} - L_{P_k} & 0 \\ \beta A P_k^T & I \end{bmatrix}^{-1} \begin{bmatrix} P_k(\beta A^T \mathbf{b} - \mathbf{g}) \\ \beta \mathbf{b} \end{bmatrix}, \end{aligned} \tag{41}$$

and hence

$$\left\{ \begin{array}{l} \mathbf{x}^{k+1} = P^k{}^T [(D_{P^k} - L_{P^k})^{-1} L_{P^k}^T] P^k \mathbf{x}^k \\ \quad + P^k{}^T (D_{P^k} - L_{P^k})^{-1} P^k (A^T \boldsymbol{\mu}^k + \beta A^T \mathbf{b} - \mathbf{g}) \\ \boldsymbol{\mu}^{k+1} := \boldsymbol{\mu}^k - \beta(A\mathbf{x}^{k+1} - \mathbf{b}). \end{array} \right. \tag{42}$$

The first equation in (42) coincides exactly with one iteration of the RSSOR with $\omega = 1$ (see Eq. (26)) for the solution of the linear system $H_\beta \mathbf{x} = A^T \boldsymbol{\mu}^k + \beta A^T \mathbf{b} - \mathbf{g}$. On the other hand, as proved in Theorem 17, the number of RSSOR sweeps per iALM-step sufficient to obtain an expected residual which ensures the a.s. convergence, is uniformly bounded above by a constant. We find that this is a noteworthy improvement of the results obtained in [12, 44, 53]. Indeed, in these works, only the *the convergence in expectation* of the iterates produced by (42) has been proved, i.e., the convergence

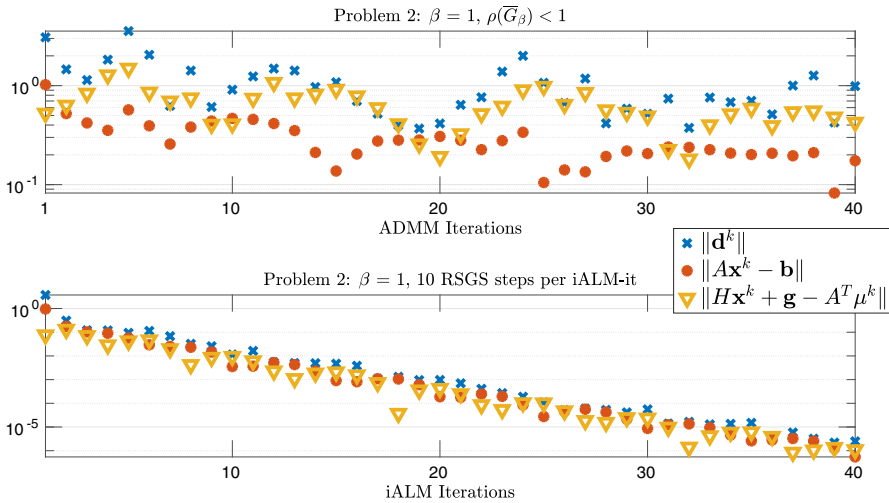


Fig. 6 Random ADMM vs iALM&RSGS for Problem 2 (logarithmic scale on y-axis)

to zero of $\|\mathbb{E}\left(\begin{bmatrix} \mathbf{x}^k \\ \boldsymbol{\mu}^k \end{bmatrix}\right) - \begin{bmatrix} \bar{\mathbf{x}} \\ \bar{\boldsymbol{\mu}} \end{bmatrix}\|$. To be precise, using the notation introduced in (41), the authors prove that $\bar{\rho}_\beta := \rho(\bar{G}_\beta) < 1$ where

$$\bar{G}_\beta := \mathbb{E}(G_\beta^P) = \frac{1}{|\mathcal{P}|} \sum_{P \in \mathcal{P}} G_\beta^P$$

and \mathcal{P} is a specific subset of all permutation matrices (\mathcal{P} is the subset of block permutation matrices with blocks of order n in [12, 53] and, in [44], \mathcal{P} is the subset of the permutation matrices obtained as $P = P_1 P_2$, where P_1 is a block permutation matrix with blocks of order n and P_2 is a permutation corresponding to a partition of d elements into n groups).

Overall, as already pointed out in [44, Sec. 2.2.4], the convergence in expectation may not be a good indicator of the robustness and the effectiveness of RADMM as there may exist problems characterized by a high $\|\mathbb{V}(G_\beta^P)\|$, where $\mathbb{V}(G_\beta^P)$ denotes the variance of the random variable G_β^P . We find that switching from a convergence in expectation to an a.s. convergence with provable expected worst case complexity as stated in Theorem 17, could be beneficial for the solution of such problems.

Even in this case, to further underpin the previous claim, we report in Fig. 6 the behaviour of $\|\mathbf{d}^k\|$, $\|\mathbf{A}\mathbf{x}^k - \mathbf{b}\|$ and $\|\mathbf{H}\mathbf{x}^k + \mathbf{g} - \mathbf{A}^T \boldsymbol{\mu}^k\|$ for RADMM and for iALM&RSGS where, at each inner iteration, 10 RSGS sweeps are performed. As it is clear from the comparison between the upper panels of Figs. 5 and 6 (and expected from the results obtained in [12, 53]), the introduction of a randomization procedure in the ADMM scheme is able to mitigate the divergence in the case of Problem 2. At the same time, analogously of what was observed in Fig. 5 for the deterministic case, the benefits of performing more than one RSGS sweep per iALM-step are evident (lower panel of Fig. 6).

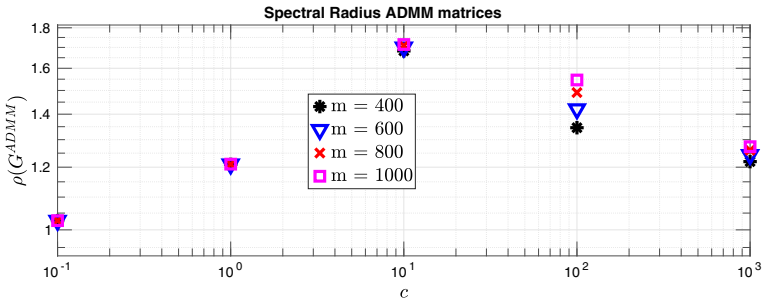


Fig. 7 Spectral radius of ADMM matrices for the generalization of Problem 2 (logarithmic scale on y-axis)

7 Numerical results

In this section we briefly present a series of numerical results aiming to guide the practitioners in the selection of some of the parameters involved in the optimal implementation of iALM&RSGS. As test set we introduce and use suitable generalizations of Problem 2. The choice of such test set is motivated from the fact that it represents, somehow, a set of pathological examples for which the convergence in expectation might not deliver satisfactory performance. To define our test set, let us introduce the matrix

$$\mathbb{R}^{d \times d} \ni \widehat{A} = \mathbf{e}\mathbf{e}^T + \begin{bmatrix} 0 & 0 & \dots & \dots & 0 \\ \vdots & \ddots & & 0 & c \\ \vdots & & 0 & \ddots & \vdots \\ 0 & c & \dots & c \end{bmatrix}.$$

We consider problem (QP) where $H = hI_{d \times d} \in \mathbb{R}^d$ with $h = 0.05$, \mathbf{b} , \mathbf{g} random vectors, and

$$A := \begin{cases} \widehat{A}, & \text{if } m = d \\ \widehat{A}(d - m + 1 : d, 1 : d), & \text{if } m \leq d \end{cases}.$$

Clearly when $m = d = 3$ and $c = 1$, we recover Problem 2. In the next Fig. 7 we plot $\rho(G^{ADMM})$ for different values of c (x -axis of the figure) and m when $d = 1000$ and $\beta = 1$ and when all the blocks are of order one (which will be precisely the setting used for the numerical results presented later).

As Fig. 7 confirms, for all the considered values of c and m we have $\rho(G^{ADMM}) > 1$, feature which endows the selected class of problems with meaningful *pathologies* suitable for testing the goodness and the robustness of our proposal when compared to RADMM.

As it is clear from Eq. (42), the dominant computational cost per RADMM step, is the solution of a block-lower triangular system performed during the GS sweep. For

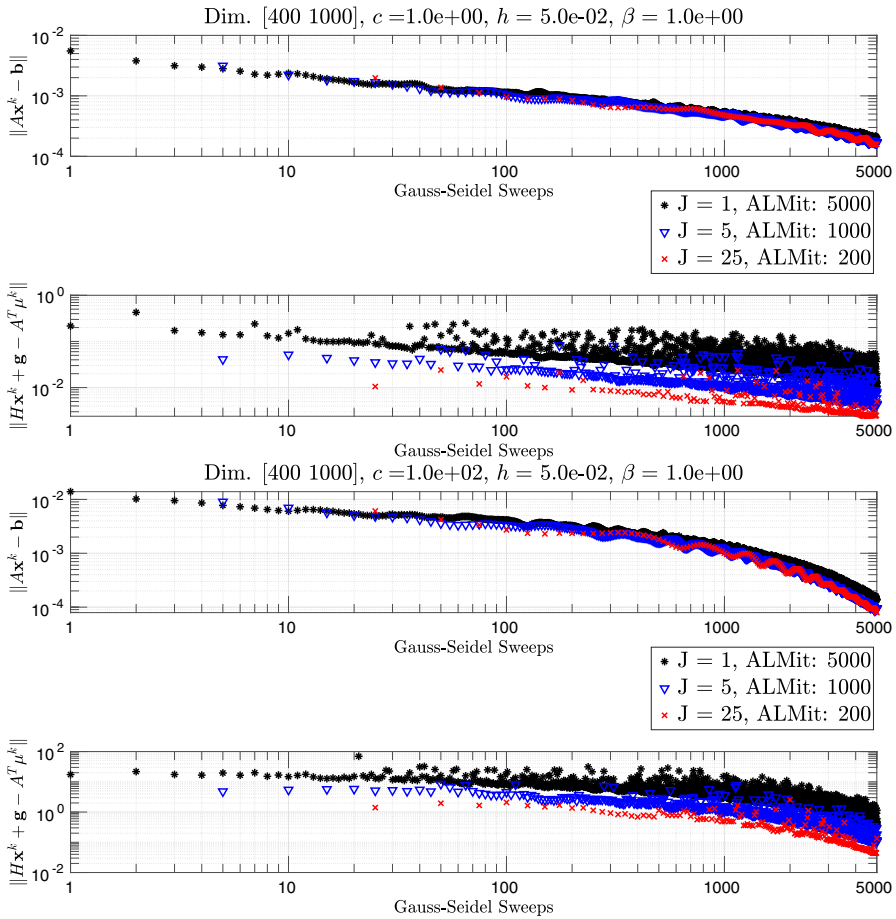


Fig. 8 Random ADMM vs iALM&RSGS(J) for selected values of m and c (logarithmic scale on both axis)

this reason, in the remainder of this section, we will fix a prescribed number of GS sweeps and measure the performance of iALM&RSGS when compared to RADMM.

In Figs. 8 and 9 we report $\|A\mathbf{x}^k - \mathbf{b}\|$ and $\|H\mathbf{x}^k + \mathbf{g} - A^T\boldsymbol{\mu}^k\|$ for RADMM and for iALM&RSGS when the maximum allowed RSGS sweeps is 5000. For the sake of brevity, we present computational results only for selected representative values of m and c ($m = 400$ and $c = 1, 100$) but a similar behaviour is observed for all the values m and c considered in Fig. 7. In particular, in Fig. 8 we present the comparison of RADMM (denoted with $J = 1$) with iALM&RSGS when $J > 1$ RSGS sweeps are performed per iALM iteration ($J = 5, 25$). In Fig. 9 instead, we compare RADMM with iALM&RSGS when RSGS for the linear system (18) is stopped if the observed residual is such that $\|\mathbf{r}^k\| < PR^{k+1}$ (see Lemma 10) with $R = 0.999$ and P is a constant depending on the initial residual $\|\mathbf{r}^0\|$.

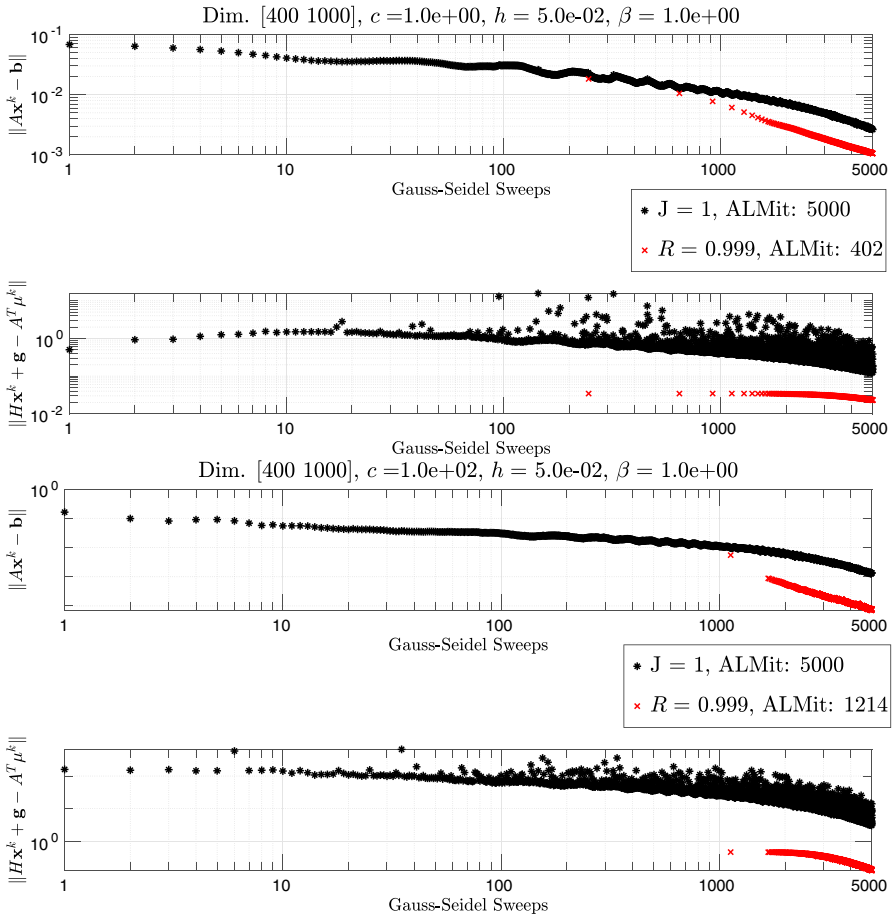


Fig. 9 RADMM vs iALM&RSGS for selected values of m and c (logarithmic scale on both axis)

Accordingly to the theoretical analysis carried out in Sect. 6, the numerical results presented in Figs. 8 and 9 confirm that performing more than one RSGS per iALM iteration consistently outperforms RADMM in the reduction of the dual residual $\|Hx^k + g - A^T \mu^k\|$. Moreover, as the comparison between Figs. 8 and 9 shows, the choice of the first strategy (fixed RSGS sweeps per iALM iteration) is preferable in general since it allows to have a faster primal/dual residuals reduction.

8 Conclusions

In this work we studied the inexact Augmented Lagrangian Method (iALM) for the solution of problem (QP). Using a splitting operator perspective, we proved that if the amount of introduced inexactness (which could be modelled also with a random variable) decreases (in expectation) accordingly to suitably chosen R^k where $R < 1$,

then we are able to give explicit asymptotic rate of convergence of the iALM (see Lemma 10). Moreover, even if the above mentioned condition requires an increasing accuracy in the linear systems to be solved at each iteration, we proved that when these linear systems are solved using the Successive-Over-Relaxation method (SOR) and its Randomly Shuffled version (RSSOR), the number of iterations sufficient to satisfy the convergence requirements can be uniformly bounded from above (see Sect. 5). Finally, using the developed theory and interpreting the n -block (Random)Alternating Direction Method of Multipliers ((R)ADMM) as an iALM which performs exactly one (RS)SOR sweep to obtain the approximate solutions of the inner linear systems, we provided computational evidence which demonstrates that the very well known convergence issues of the n -block (R)ADMM could be remedied if more than one (RS)SOR sweep for every iALM iteration were permitted (see Sect. 7).

Acknowledgements The authors are in debt with M. Rossi (University of Milano-Bicocca) for the fruitful discussions on some technical details about the probabilistic case.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

1. Arrow, K.J., Hurwicz, L., Uzawa, H.: Studies in Linear and Non-linear Programming. With contributions by H. B. Chenery, S. M. Johnson, S. Karlin, T. Marschak, R. M. Solow. Stanford Mathematical Studies in the Social Sciences, vol. II. Stanford University Press, Stanford (1958)
2. Benzi, M., Golub, G.H., Liesen, J.: Numerical solution of saddle point problems. *Acta Numer.* **14**, 1–137 (2005). <https://doi.org/10.1017/S0962492904000212>
3. Billingsley, P.: Probability and Measure. Wiley Series in Probability and Statistics. Wiley, Hoboken (2012)
4. Birken, P.: Termination criteria for inexact fixed-point schemes. *Numer. Linear Algebra Appl.* **22**(4), 702–716 (2015). <https://doi.org/10.1002/nla.1982>
5. Bodewig, E.: Matrix Calculus. Elsevier, Amsterdam (2014)
6. Boyd, S., Parikh, N., Chu, E.: Distributed Optimization and Statistical Learning Via the Alternating Direction Method of Multipliers. Now Publishers Inc., Norwell (2011)
7. Bramble, J.H., Pasciak, J.E., Vassilev, A.T.: Analysis of the inexact Uzawa algorithm for saddle point problems. *SIAM J. Numer. Anal.* **34**(3), 1072–1092 (1997). <https://doi.org/10.1137/S0036142994273343>
8. Cai, X., Han, D., Yuan, X.: On the convergence of the direct extension of ADMM for three-block separable convex minimization models with one strongly convex function. *Comput. Optim. Appl.* **66**(1), 39–73 (2017). <https://doi.org/10.1007/s10589-016-9860-y>
9. Chang, C.C., Lin, C.J.: LIBSVM: a library for support vector machines. *ACM Trans. Intell. Syst. Technol.* **2**, 27:1–27:27 (2011). Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>
10. Chen, C., Shen, Y., You, Y.: On the convergence analysis of the alternating direction method of multipliers with three blocks. *Abstr. Appl. Anal.* (2013). Art. ID 183,961, 7. <https://doi.org/10.1155/2013/183961>

11. Chen, C., He, B., Ye, Y., et al.: The direct extension of ADMM for multi-block convex minimization problems is not necessarily convergent. *Math. Program* **155**(1–2, Ser. A), 57–79 (2016). <https://doi.org/10.1007/s10107-014-0826-5>
12. Chen, C., Li, M., Liu, X., et al.: Extended ADMM and BCD for nonseparable convex minimization models with quadratic coupling terms: convergence analysis and insights. *Math. Program* **173**(1–2, Ser. A), 37–77 (2019). <https://doi.org/10.1007/s10107-017-1205-9>
13. Chen, F., Jiang, Y.L.: A generalization of the inexact parameterized Uzawa methods for saddle point problems. *Appl. Math. Comput.* **206**(2), 765–771 (2008). <https://doi.org/10.1016/j.amc.2008.09.041>
14. Chen, L., Sun, D., Toh, K.C.: An efficient inexact symmetric Gauss–Seidel based majorized ADMM for high-dimensional convex composite conic programming. *Math. Program* **161**(1–2, Ser. A), 237–270 (2017). <https://doi.org/10.1007/s10107-016-1007-5>
15. Chen, L., Li, X., Sun, D., et al.: On the equivalence of inexact proximal ALM and ADMM for a class of convex composite programming. *Math. Program* **185**(1–2, Ser. A), 111–161 (2021). <https://doi.org/10.1007/s10107-019-01423-x>
16. Chen, X.: On preconditioned Uzawa methods and SOR methods for saddle-point problems. *J. Comput. Appl. Math.* **100**(2), 207–224 (1998). [https://doi.org/10.1016/S0377-0427\(98\)00197-6](https://doi.org/10.1016/S0377-0427(98)00197-6)
17. Cheng, X.L.: On the nonlinear inexact Uzawa algorithm for saddle-point problems. *SIAM J. Numer. Anal.* **37**(6), 1930–1934 (2000). <https://doi.org/10.1137/S0036142998349266>
18. Cui, M.: A sufficient condition for the convergence of the inexact Uzawa algorithm for saddle point problems. *J. Comput. Appl. Math.* **139**(2), 189–196 (2002). [https://doi.org/10.1016/S0377-0427\(01\)00430-7](https://doi.org/10.1016/S0377-0427(01)00430-7)
19. Davis, D., Yin, W.: A three-operator splitting scheme and its optimization applications. *Set-Valued Var. Anal.* **25**(4), 829–858 (2017). <https://doi.org/10.1007/s11228-017-0421-z>
20. Durrett, R.: *Probability: Theory and Examples*. Cambridge Series in Statistical and Probabilistic Mathematics, vol. 31, 4th edn. Cambridge University Press, Cambridge (2010). <https://doi.org/10.1017/CBO9780511779398>
21. Eckstein, J.: *Splitting methods for monotone operators with applications to parallel optimization*. PhD thesis, Massachusetts Institute of Technology (1989)
22. Eckstein, J., Yao, W.: Augmented Lagrangian and alternating direction methods for convex optimization: a tutorial and some illustrative computational results. *RUTCOR Res. Rep.* **32**(3), 44 (2012)
23. Eckstein, J., Yao, W.: Understanding the convergence of the alternating direction method of multipliers: theoretical and computational perspectives. *Pac. J. Optim.* **11**(4), 619–644 (2015)
24. Elman, H.C., Golub, G.H.: Inexact and preconditioned Uzawa algorithms for saddle point problems. *SIAM J. Numer. Anal.* **31**(6), 1645–1661 (1994). <https://doi.org/10.1137/0731085>
25. Fortin, M., Glowinski, R.: *Augmented Lagrangian Methods: Applications to the Numerical Solution of Boundary-Value Problems*. Elsevier, Amsterdam (2000)
26. Frankel, S.P.: Convergence rates of iterative treatments of partial differential equations. *Math. Tables Aids Comput.* **4**, 65–75 (1950)
27. Gauss, C.F.: *Werke* (in German), vol. 9. Königlich-Gelehrten Gesellschaft der Wissenschaften, Göttingen, pp. 763–764 (1903)
28. Glowinski, R., Marrocco, A.: Sur l’approximation, par éléments finis d’ordre un, et la résolution, par pénalisation-dualité, d’une classe de problèmes de Dirichlet non linéaires. *Rev Française Automat Informat Recherche Opérationnelle Sér Rouge Anal Numér* **9**(R-2), 41–76 (1975)
29. Goldfarb, D., Ma, S.: Fast multiple-splitting algorithms for convex optimization. *SIAM J. Optim.* **22**(2), 533–556 (2012). <https://doi.org/10.1137/090780705>
30. Golub, G.H., Wu, X., Yuan, J.Y.: SOR-like methods for augmented systems. *BIT* **41**(1), 71–85 (2001). <https://doi.org/10.1023/A:1021965717530>
31. Hackbusch, W.: *Iterative Solution of Large Sparse Systems of Equations*. Applied Mathematical Sciences, vol. 95, 2nd edn. Springer, Berlin (2016). <https://doi.org/10.1007/978-3-319-28483-5>
32. Hager, W.W., Zhang, H.: Convergence rates for an inexact ADMM applied to separable convex optimization. *Comput. Optim. Appl.* **77**(3), 729–754 (2020). <https://doi.org/10.1007/s10589-020-00221-y>
33. Han, D., Yuan, X.: A note on the alternating direction method of multipliers. *J. Optim. Theory Appl.* **155**(1), 227–238 (2012). <https://doi.org/10.1007/s10957-012-0003-z>
34. He, B., Tao, M., Yuan, X.: Alternating direction method with Gaussian back substitution for separable convex programming. *SIAM J. Optim.* **22**(2), 313–340 (2012). <https://doi.org/10.1137/110822347>

35. He, B., Tao, M., Xu, M., et al.: An alternating direction-based contraction method for linearly constrained separable convex programming problems. *Optimization* **62**(4), 573–596 (2013). <https://doi.org/10.1080/02331934.2011.611885>
36. Hestenes, M.R.: Multiplier and gradient methods. *J. Optim. Theory Appl.* **4**, 303–320 (1969). <https://doi.org/10.1007/BF00927673>
37. Horn, R.A., Johnson, C.R.: *Matrix Analysis*, 2nd edn. Cambridge University Press, Cambridge (2013)
38. Kang, M., Kang, M., Jung, M.: Inexact accelerated augmented Lagrangian methods. *Comput. Optim. Appl.* **62**(2), 373–404 (2015). <https://doi.org/10.1007/s10589-015-9742-8>
39. Lan, G., Monteiro, R.D.C.: Iteration-complexity of first-order augmented Lagrangian methods for convex programming. *Math. Program* **155**(1–2, Ser. A), 511–547 (2016). <https://doi.org/10.1007/s10107-015-0861-x>
40. Li, J., Hu, Z.C.: Toeplitz lemma, complete convergence, and complete moment convergence. *Commun. Stat. Theory Methods* **46**(4), 1731–1743 (2017). <https://doi.org/10.1080/03610926.2015.1026996>
41. Li, M., Sun, D., Toh, K.C.: A convergent 3-block semi-proximal ADMM for convex minimization problems with one strongly convex block. *Asia-Pac. J. Oper. Res.* **32**(4), 1550,024 (2015). <https://doi.org/10.1142/S0217595915500244>
42. Lin, T., Ma, S., Zhang, S.: On the global linear convergence of the ADMM with multiblock variables. *SIAM J. Optim.* **25**(3), 1478–1497 (2015). <https://doi.org/10.1137/140971178>
43. Liu, Y.F., Liu, X., Ma, S.: On the nonergodic convergence rate of an inexact augmented Lagrangian framework for composite convex programming. *Math. Oper. Res.* **44**(2), 632–650 (2019). <https://doi.org/10.1287/moor.2018.0939>
44. Mihic, K., Zhu, M., Ye, Y.: Managing randomization in the multi-block alternating direction method of multipliers for quadratic optimization. *Math. Program Comput.* (2020). <https://doi.org/10.1007/s12532-020-00192-5>
45. Nedelcu, V., Necoara, I., Tran-Dinh, Q.: Computational complexity of inexact gradient augmented Lagrangian methods: application to constrained MPC. *SIAM J. Control Optim.* **52**(5), 3109–3134 (2014). <https://doi.org/10.1137/120897547>
46. Ortega, J.M., Rheinboldt, W.C.: *Iterative Solution of Nonlinear Equations in Several Variables*. Academic Press, New York (1970)
47. Oswald, P.: On the convergence rate of SOR: a worst case estimate. *Computing* **52**(3), 245–255 (1994). <https://doi.org/10.1007/BF02246506>
48. Oswald, P., Zhou, W.: Random reordering in SOR-type methods. *Numer. Math.* **135**(4), 1207–1220 (2017). <https://doi.org/10.1007/s00211-016-0829-7>
49. Peng, Y., Ganesh, A., Wright, J., et al.: RASL: robust alignment by sparse and low-rank decomposition for linearly correlated images. *IEEE Trans. Pattern Anal. Mach. Intell.* **34**(11), 2233–2246 (2012)
50. Powell, M.J.D.: A Method for Nonlinear Constraints in Minimization Problems. In: Fletcher, R., (Ed.) *Optimization*. Academic Press, New York, NY, pp. 283–298 (1969)
51. Ryu, E.K., Boyd, S.: A primer on monotone operator methods (survey). *Appl. Comput. Math.* **15**(1), 3–43 (2016)
52. Stout, W.F.: *Almost Sure Convergence, Probability and Mathematical Statistics*, vol. 24. Academic Press [A Subsidiary of Harcourt Brace Jovanovich, Publishers], New York (1974)
53. Sun, R., Luo, Z.Q., Ye, Y.: On the efficiency of random permutation for ADMM and coordinate descent. *Math. Oper. Res.* **45**(1), 233–271 (2020). <https://doi.org/10.1287/moor.2019.0990>
54. Tao, M.: Convergence study of indefinite proximal ADMM with a relaxation factor. *Comput. Optim. Appl.* **77**(1), 91–123 (2020). <https://doi.org/10.1007/s10589-020-00206-x>
55. Tao, M., Yuan, X.: Recovering low-rank and sparse components of matrices from incomplete and noisy observations. *SIAM J. Optim.* **21**(1), 57–81 (2011). <https://doi.org/10.1137/100781894>
56. Varga, R.S.: *Matrix Iterative Analysis*. Prentice-Hall Inc., Englewood Cliffs (1962)
57. Xu, Y.: Iteration complexity of inexact augmented Lagrangian methods for constrained convex programming. *Math. Program* **185**(1–2, Ser. A), 199–244 (2021). <https://doi.org/10.1007/s10107-019-01425-9>
58. Young, D.M.: *Iterative methods for solving partial difference equation of elliptic type*. ProQuest LLC, Ann Arbor, MI, thesis (Ph.D.)—Harvard University (1950)