

A parallelizable augmented Lagrangian method applied to large-scale non-convex-constrained optimization problems

Natashia Boland¹ · Jeffrey Christiansen² · Brian Dandurand² · Andrew Eberhard² · Fabricio Oliveira³

Received: 20 November 2016 / Accepted: 22 February 2018 / Published online: 1 March 2018
© Springer-Verlag GmbH Germany, part of Springer Nature and Mathematical Optimization Society 2018

Abstract We contribute improvements to a Lagrangian dual solution approach applied to large-scale optimization problems whose objective functions are convex, continuously differentiable and possibly nonlinear, while the non-relaxed constraint set is compact but not necessarily convex. Such problems arise, for example, in the split-variable deterministic reformulation of stochastic mixed-integer optimization problems. We adapt the augmented Lagrangian method framework to address the presence of nonconvexity in the non-relaxed constraint set and to enable efficient parallelization. The development of our approach is most naturally compared with the development of proximal bundle methods and especially with their use of serious step conditions. However, deviations from these developments allow for an improvement in efficiency with which parallelization can be utilized. Pivotal in our modification to the augmented Lagrangian method is an integration of the simplicial decomposition method and the nonlinear block Gauss–Seidel method. An adaptation of a serious step condition associated with proximal bundle methods allows for the approximation tolerance to be automatically adjusted. Under mild conditions optimal dual convergence is proven, and we report computational results on test instances from the stochas-

This work was supported by the Australian Research Council (ARC) Grant ARC DP140100985.

Electronic supplementary material The online version of this article (<https://doi.org/10.1007/s10107-018-1253-9>) contains supplementary material, which is available to authorized users.

✉ Andrew Eberhard
andy.eberhard@rmit.edu.au

¹ Georgia Institute of Technology, Atlanta, GA, USA

² RMIT University, Melbourne, VIC, Australia

³ Aalto University, Espoo, Finland

tic optimization literature. We demonstrate improvement in parallel speedup over a baseline parallel approach.

Keywords Augmented Lagrangian method · Proximal bundle method · Nonlinear block Gauss–Seidel method · Simplicial decomposition method · Parallel computing

Mathematics Subject Classification 90-08 · 90C06 · 90C11 · 90C15 · 90C25 · 90C26 · 90C30 · 90C46

1 Introduction

We develop a dual solution approach to the problem of interest having the form

$$\zeta^* := \min_{x,z} \{f(x) : Qx = z, x \in X, z \in Z\}, \quad (1)$$

where f is convex and continuously differentiable, $Q \in \mathbb{R}^{q \times n}$ is a block-diagonal matrix determining linear constraints $Qx = z$, $X \subset \mathbb{R}^n$ is a closed and bounded set, and $Z \subset \mathbb{R}^q$ is a linear subspace. The vector $x \in X$ of decision variables is derived from the original decisions associated with a problem, while the vector $z \in Z$ of auxiliary variables are introduced to effect a decomposable structure in (1). In particular, the decomposable structure takes the form: (1) $X = \prod_{i=1}^m X_i$ with $X_i \subset \mathbb{R}^{n_i}$ closed and bounded and $\sum_{i=1}^m n_i = n$; (2) $f(x) = \sum_{i=1}^m f_i(x_i)$ where $f_i : \mathbb{R}^{n_i} \mapsto \mathbb{R}$ are convex and differentiable for $i = 1, \dots, m$; (3) Q has block diagonal structure with block diagonal components denoted as $Q_i \in \mathbb{R}^{q_i \times n_i}$, $i = 1, \dots, m$ where $\sum_{i=1}^m q_i = q$, so that after setting $z = (z_i)_{i=1, \dots, m}$, where for each $i = 1, \dots, m$, $z_i \in \mathbb{R}^{q_i}$, we may write $Qx = z$ as $Q_i x_i = z_i$, $i = 1, \dots, m$. This decomposable structure is implicitly present throughout the development of this paper, although explicit referral to it is typically avoided where it is not needed. We make no other explicit assumptions on f , Q , X , or Z , but we otherwise assume that problem (1) is feasible with finite optimal value.

Problem (1) is general enough to subsume the split-variable deterministic reformulation of a stochastic optimization problem with potentially multiple stages, as defined, for example, in [1], while it can also model the case where f is nonlinear and convex and/or X is any compact (but not necessarily convex) set.

We denote the convex hull of X by $\text{conv}(X)$. The nonconvexity of X is avoided in our development in the sense that (1) an oracle is assumed to be provided for solving a subproblem with linear objective over feasible set X , and (2) another oracle is assumed to be provided for solving a subproblem with convex (possibly nonlinear) objective over a closed, convex, and nonempty subset of $\text{conv}(X)$.

We develop a solution approach to solving the following relaxation of (1),

$$\zeta^{CLD} := \min_{x,z} \{f(x) : Qx = z, x \in \text{conv}(X), z \in Z\} \quad (2)$$

and its Lagrangian dual problem due to the relaxation of $Qx = z$,

$$\zeta^{CLD} = \sup_{\omega} \phi^C(\omega), \quad (3)$$

which is based on the dual function

$$\phi^C(\omega) := \min_x \left\{ f(x) + \omega^\top (Qx - z) : x \in \text{conv}(X), z \in Z \right\}. \quad (4)$$

When f is linear, then $\phi^C(\omega) = \phi(\omega)$ where

$$\phi(\omega) := \min_{x,z} \left\{ f(x) + \omega^\top (Qx - z), x \in X, z \in Z \right\}. \quad (5)$$

(That is, when f is linear, the role of X and $\text{conv}(X)$ are interchangeable.) Consequently, when f is linear, $\zeta^{CLD} = \zeta^{LD} := \sup_{\omega} \phi(\omega)$. However, when f is nonlinear, then in general, $\zeta^{CLD} \leq \zeta^{LD}$. Strict inequality is demonstrated with the following example. Let $f : \mathbb{R}^2 \mapsto \mathbb{R}$ be defined by $f(x) = (x_1 - 0.5)^2 + (x_2 - 0.5)^2$, $X = \{0, 1\} \times \{0, 1\}$, and let $Qx = z$ be defined to model the constraints $x_1 - z_1 = 0$ and $x_2 - z_2 = 0$ where $Z = \{(z_1, z_2) : z_1 = z_2\} \subset \mathbb{R}^2$. We see trivially that $\zeta^{CLD} = 0$, which is verified with the saddle point $x_1^* = x_2^* = z_1^* = z_2^* = 0.5$ and $\omega^* = (0, 0)$. However, $\zeta^{LD} = 0.5$, which is verified with either of the saddle points $x_1^* = x_2^* = z_1^* = z_2^* = 0$ and $\omega^* = (0, 0)$, or $x_1^* = x_2^* = z_1^* = z_2^* = 1$ and $\omega^* = (0, 0)$. Thus, $\zeta^{CLD} < \zeta^{LD}$.

Given that X is compact and f is continuous, in order for $-\infty < \phi^C(\omega)$ to hold, it is necessary and sufficient that the dual feasibility assumption

$$\omega \in Z^\perp := \left\{ v \in \mathbb{R}^q : v^\top z = 0 \text{ for all } z \in Z \right\} \quad (6)$$

is maintained either by assumption or by construction. Under condition (6), the z term in definition (4) vanishes, and we may compute

$$\phi^C(\omega) = \min_x \left\{ f(x) + \omega^\top Qx : x \in \text{conv}(X) \right\}.$$

Consequently, ϕ^C becomes separable as

$$\phi^C(\omega) = \sum_{i=1}^m \phi_i^C(\omega_i),$$

where $\phi_i^C(\omega_i) := \min_x \left\{ f_i(x_i) + \omega_i^\top Q_i x_i : x_i \in \text{conv}(X_i) \right\}$ and $\omega = (\omega_1, \dots, \omega_m) \in \mathbb{R}^{q_1} \times \dots \times \mathbb{R}^{q_m}$ has a block structure compatible with the block diagonal structure of Q .

Given that X is closed and bounded [thus so is $\text{conv}(X)$], and (2) is assumed to be feasible, then in order to guarantee that the maximum in (3) is realized for some $\omega^* \in Z^\perp$, we assume a constraint qualification such as Slater's condition. In

other words, we assume that there exists (x^*, z^*) such that $x^* \in \text{int}(\text{conv}(X))$ and $Qx^* = z^*$, where $\text{int}(\cdot)$ returns the topological interior of the set argument. If $\text{conv}(X)$ is polyhedral, then even this Slater's condition is not required.

For the Lagrangian dual problem (3), we note that the objective function ϕ^C is concave, even when f is not convex. We can apply a subgradient method (see e.g. [2]; in textbooks [3, 4]) for solving (3) in an efficiently parallelizable manner. Such an approach is proposed in [5]. However, it is preferable to make use of structural features of (3) that allow for smoothing or regularization, so that better convergence properties are realized. For this reason, we consider alternative developments based on proximal point methods that are modified to address both of the above two challenges.

In this paper, we develop an iterative solution approach to solving problems (2) and (3) subject to the two main challenges. First, the set X is not convex (for example, it may have mixed-integer constraints as part of its definition), $\text{conv}(X)$ is not given explicitly, and thus, the augmented Lagrangian method is not theoretically supported with constraint set X , or implementable¹ with constraint set $\text{conv}(X)$. Second, the solution approach should be amenable to efficient parallel computation, in the sense of maximizing the computational work that can be parallelized, the memory usage that can be distributed, and minimizing the amount of parallel communication.

The remainder of the paper is organized as follows. In Sect. 2, we provide a background and literature review. In Sect. 3, a general algorithmic framework based on the AL method with approximate subproblem solutions is developed and analyzed. In Sect. 4, a specific implementation of the Sect. 3 framework is posed based on the integration of SDM and GS methods, which addresses the aforementioned issues of implementability and efficiency of parallelization. In Sect. 5, computational experiments and their outcomes are described and interpreted. And at last, Sect. 6 concludes the paper and provides avenues for future work.

2 Background

As a starting point, we first consider the classical augmented Lagrangian method based on proximal point methods. The augmented Lagrangian (AL) method (also known as the method of multipliers) is developed from proximal point methods, and references include [3, 6–8].

The AL method typically has favorable convergence properties as a dual solution approach for convex problems (linear convergence rate under certain assumptions, see [6, 9] and references cited therein). However, two issues arise in the setting we consider: (1) the set X is not convex, and so current theories of convergence are not applicable; and (2) the primal subproblem associated with each iteration of the AL method is not separable due to the augmented Lagrange term, making efficient parallel implementations difficult to develop.

The proximal bundle method initially appeared in [10], and for a survey, which also gives its history, see [11]. Use of inexact oracles for computing $\phi(\omega)$ and elements

¹ Implementable is used in the sense of ability to implement the required functionality of an algorithm, not the meaning of implementable specific to stochastic optimization.

of the subdifferential set $\partial\phi(\omega)$ are studied in [11–13] and references therein. In its dual form, the bundle method may be referred to as the stabilized column generation method [14] or the proximal simplicial decomposition method [15]. In implementation, the algorithm we develop here more closely resembles the latter dual form.

For parallelization of the proximal bundle method, see [16] and [17]. The approach developed in this paper is most naturally compared with [17], as both approaches address the manner in which the same continuous master problem is approximately solved. The approach of [16] uses a substantially different parallel computational paradigm based on subspace optimization. This approach, in which solution subspaces are assigned to processors based on periodically updated global state information, is not necessarily based on the problem's decomposable structure.

The proximal bundle method approach requires modification for efficient parallelization. This matter is addressed in [17], where a solution to the continuous master problem is obtained by primal dual interior point methods that exploit the decomposable structure present in the augmented Lagrangian term. We provide and analyze an alternative approach based on the use of:

1. the simplicial decomposition method (SDM) [3, 18–20], which provides an alternative framework to the proximal bundle method to address the implementability of the proximal point method while allowing for the possibility that f is nonlinear; and
2. nonlinear block Gauss–Seidel (GS) method [21–25] to approximate the solutions to the continuous master problem.

Motivated by its constituent parts, the algorithm we develop is referred to as SDM-GS-ALM.

In an iteration of SDM-GS-ALM, the analog to the continuous master problem is not solved to (near) exactness; instead, approximate solutions based on possibly just one nonlinear block GS iteration are used. Due to the underlying need for convexification of the non-relaxed constraint set, implementability requires that the nonlinear block GS method must be integrated with the SDM so that optimal convergence of the resulting iterations can be established. In this way, a serious step condition similar to that used in proximal bundle methods is eventually satisfied after a finite number of such integrated SDM-GS iterations, and analogous dual optimal convergence of our approach is recovered even with the deviations from the proximal bundle method. In summary, we algorithmically integrate the AL method, the SDM, nonlinear block GS iterations, and the proximal bundle method serious step condition. A convergence analysis is also provided for SDM-GS-ALM. Such an integration allows for a considerable improvement in parallel efficiency with respect to maximizing the computational work that can be parallelized, the memory usage that can be distributed, and minimizing the amount of parallel communication.

Other methods developed in the past that are related to aspects of our contribution include the following. In terms of approximating within the AL method, we include reference to [26, 27], where the research goal of developing implementable approximation criteria is addressed. The separable augmented Lagrangian (SALA) method [28], which is an application of the alternating direction method of multipliers

(ADMM) [29–31] with a form of resource allocation decomposition and incorporates separability into the AL method. Other approaches to introducing separability into the AL method include [32, 33]. Jacobi iterate approaches applied within either a proximal bundle method or an AL method framework are considered in [34, 35]; the accelerated distributed augmented Lagrangian method (ADAL) developed in [32] is like a Jacobi-iterate analogue of ADMM with supporting convergence analysis. Other approaches to incorporating separability are found in the alternating linearization approaches [36, 37] and the predictor corrector proximal multiplier (PCPM) methods [38, 39]. All of these methods provide implementable mechanisms for approximating primal subproblem solutions and effecting parallelism in a setting where X is convex. However, they are not practically implementable in our setting where X is not convex and its convex hull $\text{conv}(X)$ is not given beforehand in a computationally useful closed-form description.

Another recently developed algorithm, referred to as FW-PH [40], is closely related to the SDM-GS-ALM algorithm developed in this paper. In terms of functionality, both appear as modifications to ADMM with inner approximated subproblem solutions. While the algorithms differ only slightly in terms of functionality, there are substantial differences in the motivation and the convergence analysis. The convergence analysis of FW-PH interfaces with the convergence analysis for ADMM, which is most naturally developed in the context of the theory of maximal monotone operators and Douglas–Rachford splitting methods [41, 42], or as the proximal decomposition of the graph of a maximal monotone operator [43]. In contrast, the convergence analysis of SDM-GS-ALM naturally reflects its synthesis of SDM, the nonlinear block GS method, the proximal bundle method, and the AL method. The convergence analysis of SDM-GS-ALM follows under more general assumptions than that for FW-PH. In particular, the convergence analysis of SDM-GS-ALM allows for trimming of the inner approximations, and it does not require the warm-starting required by FW-PH. The most important difference in functionality is due to the influence of ideas from proximal bundle methods in SDM-GS-ALM, where updates of ω are taken conditionally at each iteration, while such updates are taken unconditionally at each iteration of FW-PH. We shall see that these conditional updates help to mitigate performance problems that arise due to the seemingly inevitable use of suboptimal algorithm parameters.

In papers such as [44, 45], ADMM is applied directly to the primal problem (1). In both works, it is acknowledged that ADMM is not theoretically supported in optimal convergence due to the lack of convexity of X . Nevertheless, [45] reports the potential for Lagrangian dual bounds to be recovered at each iteration of ADMM even though it is applied to (1). In [44], where ADMM is applied to nonconvex decentralized unit commitment problems, heuristic improvements to ADMM are introduced to address the lack of convexity due to the mixed-integer constraints. In contrast to both of these approaches, where ADMM is applied directly to the primal problem (1), the approach developed in this paper, and its related approach [40], both resemble ADMM

but with application to a primal characterization of the dual problem. In these two approaches, the challenge of not having an explicit form for this primal characterization is addressed.

3 An alternative AL approximation approach

Algorithm 1 provides a general framework for an AL method with approximate subproblem solutions that uses the bundle method's serious step condition (SSC). This framework will be useful to guide the developments presented next, in which we discuss optimality conditions and convergence properties of the algorithm. The convergence proof of Algorithm 1 is based on the convergence proofs of the proximal bundle method such as found in Chapter 7 of [4]. In the following we denote the augmented Lagrangian (AL):

$$\begin{aligned} L_\rho(x, z, \omega) &:= f(x) + \omega^\top(Qx - z) + \frac{\rho}{2}\|Qx - z\|^2 \\ &= f(x) + \omega^\top Qx + \frac{\rho}{2}\|Qx - z\|^2 \end{aligned} \quad (7)$$

where the AL relaxes $Qx = z$ and with the second equality following from $\omega \in Z^\perp$ and $z \in Z$.

In the proximal bundle method, the dual function ϕ is approximated by a cutting plane model function that majorizes ϕ . Instead we use the following approximation $\hat{\phi} : \mathbb{R}^q \times \mathbb{R}^n \times \mathbb{R}^q \mapsto \mathbb{R}$ of ϕ^C centered at (x^k, z^k) , $k \geq 0$, to replace the cutting plane model:

$$\hat{\phi}(\omega, x^k, z^k) := L_\rho(x^k, z^k, \omega) + \frac{\rho}{2}\|Qx^k - z^k\|_2^2.$$

The convex hull $\text{conv}(X)$ is not known explicitly, and so ϕ^C cannot be evaluated directly. Consequently, we additionally make use of the following minorization $\check{\phi}$ of ϕ^C that can be evaluated. For $x^k \in \text{conv}(X)$, $k \geq 0$, define $\check{\phi}(\omega, x^k)$ as follows:

$$\check{\phi}(\omega, x^k) := \min_x \left\{ f(x^k) + \nabla_x f(x^k)^\top(x - x^k) + \omega^\top Qx : x \in X \right\}. \quad (8)$$

In Algorithm 1, a proximal bundle method-like serious step condition is used in Line 9 that makes use of $\hat{\phi}$ and $\check{\phi}$ in place of the cutting plane model and ϕ , respectively. The inputs f , Q , X , and Z specify the data associated with problem (1); $\rho > 0$ is the AL term coefficient; ω^0 is an initial dual solution; $\gamma \in (0, 1)$ is the parameter of the serious step condition of Line 9; and $\epsilon > 0$ is a tolerance for termination. Algorithm 1 will be given a specific implementation in the form of SDM-GS-ALM in Sect. 4.

Algorithm 1 A general approximated ALM using a bundle method SSC.

```

1: Preconditions:  $\omega^1 \in Z^\perp$ ,  $\gamma \in (0, 1)$ .
2: function APPROXALM( $f, Q, X, Z, \rho, \omega^1, \gamma, \epsilon, k_{\max}$ )
3:   for  $k = 1, 2, \dots, k_{\max}$  do
4:     Solve approximately
5:      $(x^k, z^k) \in \operatorname{argmin}_{x,z} \left\{ L_\rho(x, z, \omega^k) : x \in \operatorname{conv}(X), z \in Z \right\}$  such that
6:       1)  $z^k \in \operatorname{argmin}_z \left\{ \|Qx^k - z\|_2^2 : z \in Z \right\}$  and
7:       2) either
8:          $\widehat{\phi}(\omega^k, x^k, z^k) - \check{\phi}(\omega^k, x^{k-1}) \leq \epsilon$  or
9:          $0 < \gamma \leq \frac{\check{\phi}(\omega^k + \rho(Qx^k - z^k), x^k) - \check{\phi}(\omega^k, x^{k-1})}{\widehat{\phi}(\omega^k, x^k, z^k) - \check{\phi}(\omega^k, x^{k-1})}$ 
10:      if  $\widehat{\phi}(\omega^k, x^k, z^k) - \check{\phi}(\omega^k, x^{k-1}) \leq \epsilon$  then
11:        return  $(x^k, z^k, \omega^k)$ 
12:      else
13:        set  $\omega^{k+1} \leftarrow \omega^k + \rho(Qx^k - z^k)$ 
14:      end if
15:    end for
16:  return  $(x^k, z^k, \omega^{k+1})$ 
17: end function

```

Remark 1 Under the assumption that $\operatorname{conv}(X)$ is not known beforehand by any characterization, direct evaluation of ϕ^C or any of its subgradients at any $\omega \in Z^\perp$ is not possible. This dual function is not used in the proximal bundle method and is only treated indirectly in the current development.

In addition to generating a sequence $\{\omega^k\}$ of dual solutions to (3), our algorithm will also generate a sequence of primal solutions $\{(x^k, z^k)\}$ to (2), and so reference to (2) will be useful. In applying the AL method to problem (2), the continuous master problem for fixed $\omega \in Z^\perp$ takes the form

$$\zeta_\rho^{AL}(\omega) := \min_{x,z} \left\{ L_\rho(x, z, \omega), x \in \operatorname{conv}(X), z \in Z \right\}. \quad (9)$$

Lemma 1 For any optimal solution ω^* to problem (3), we have $\zeta_\rho^{AL}(\omega^*) = \zeta^{CLD}$. Additionally, any optimal solution (x^*, z^*) to problem (9) with $\omega = \omega^*$ is also optimal for problem (2).

Proof We specialize developments in, e.g., Section 4 of [46] or Section 6.4.3 of [4]. In the following, we begin assuming $\bar{\omega} \in Z^\perp$ as an arbitrary fixed vector to show:

$$\max_{\omega \in Z^\perp} \left\{ (\omega - \bar{\omega})^\top Qx - \frac{1}{2\rho} \|\omega - \bar{\omega}\|_2^2 \right\} = \frac{\rho}{2} \min_{z \in Z} \|Qx - z\|_2^2. \quad (10)$$

By the uniqueness of the projection onto a subspace, $z \in \operatorname{argmin}_{z \in Z} \|Qx - z\|_2^2$ is the unique $z \in Z$ for which $(Qx - z) \in Z^\perp$. Moreover, the optimality condition for left-hand side problem in (10): $0 \in \partial_\omega \left\{ (\omega - \bar{\omega})^\top Qx - \frac{1}{2\rho} \|\omega - \bar{\omega}\|_2^2 + \delta_{Z^\perp}(\omega) \right\}$ dictate that $(\omega - \bar{\omega}) \in \rho(Qx + Z)$ (where we have used, for $\omega \in Z^\perp$, $\partial \delta_{Z^\perp}(\omega) = N_{Z^\perp}(\omega) =$

Z , where δ_{Z^\perp} denotes the indicator function of a convex set Z^\perp and $N_{Z^\perp}(\omega)$ the normal cone at ω). Hence $(\omega - \bar{\omega}) = \rho(Qx - z)$ for some $z \in Z$. Furthermore, $\omega - \bar{\omega} \in Z^\perp$ must hold, and so $z \in Z$ must be chosen so that $\rho(Qx - z) \in Z^\perp$ also. Consequently, from our first observation, this $z \in Z$ must be the unique solution of the right-hand side of (10). Evaluating the objective on the left-hand side of (10) (observing that we must have $\rho(Qx - z) \in Z^\perp$) establishes the claimed equality. We may compute:

$$\begin{aligned}
 & \max_{\omega \in Z^\perp} \phi^C(\omega) - \frac{1}{2\rho} \|\omega - \bar{\omega}\|_2^2 \\
 &= \max_{\omega \in Z^\perp} \min_x \left\{ f(x) + \omega^\top Qx - \frac{1}{2\rho} \|\omega - \bar{\omega}\|_2^2 : x \in \text{conv}(X) \right\} \\
 &= \min_x \left\{ f(x) + \bar{\omega}^\top Qx + \max_{\omega \in Z^\perp} \left\{ (\omega - \bar{\omega})^\top Qx - \frac{1}{2\rho} \|\omega - \bar{\omega}\|_2^2 \right\} : x \in \text{conv}(X) \right\} \\
 &= \min_x \left\{ f(x) + \bar{\omega}^\top Qx + \frac{\rho}{2} \min_{z \in Z} \left\{ \|Qx - z\|_2^2 \right\} : x \in \text{conv}(X) \right\} \\
 &= \min_{x, z} \left\{ L_\rho(x, z, \bar{\omega}), x \in \text{conv}(X), z \in Z \right\}. \tag{11}
 \end{aligned}$$

The switching of min and max is justified by the Sion min–max theorem along with the convexity of f , $\text{conv}(X)$, Z and $\text{conv}(X)$ assumed compactness. In substituting $\bar{\omega} = \omega^*$, the value of the left-hand side maximization problem (11) is clearly ζ^{CLD} , while the same substitution on the right-hand side (10) yields the value $\zeta_\rho^{AL}(\omega^*)$, from which we see that $\zeta^{CLD} = \zeta_\rho^{AL}(\omega^*)$. To prove the last claim, we note that $L_\rho(x^*, z^*, \omega^*) = \zeta^{CLD}$ implies that $\|Qx^* - z^*\|_2^2 = 0$. Otherwise, $\phi^C(\omega^*) < \zeta^{CLD}$, contradicting the dual optimality of ω^* . Thus, (x^*, z^*) is feasible and optimal for problem (2). \square

The approximation $\hat{\phi}$ satisfies the following bounding relationship.

Lemma 2 For each (x^k, z^k) , $k \geq 0$, such that the z -optimality condition is satisfied:

$$z^k \in \operatorname{argmin}_z \left\{ \|Qx^k - z\|_2^2 : z \in Z \right\}, \tag{12}$$

we have for each $\omega \in Z^\perp$

$$\hat{\phi}(\omega, x^k, z^k) \geq \phi^C \left(\omega + \rho \left(Qx^k - z^k \right) \right). \tag{13}$$

Proof Via convexity of the term $\|Qx - z\|_2^2$ over $(x, z) \in \text{conv}(X) \times Z$, we may write the following inequalities that hold for $(x, z) \in \text{conv}(X) \times Z$ and a fixed $\omega \in Z^\perp$. Via the subgradient inequality:

$$\begin{aligned}
L_\rho(x, z, \omega) &\geq f(x) + \omega^\top Qx + \frac{\rho}{2} \|Qx^k - z^k\|_2^2 \\
&\quad + \rho \left(Qx^k - z^k \right)^\top (Qx - z) - \rho \left(Qx^k - z^k \right)^\top (Qx^k - z^k) \\
&= f(x) + \omega^\top Qx - \frac{\rho}{2} \|Qx^k - z^k\|_2^2 + \rho \left(Qx^k - z^k \right)^\top (Qx - z) \\
&\implies L_\rho(x, z, \omega) + \frac{\rho}{2} \|Qx^k - z^k\|_2^2 \geq f(x) + \left[\omega + \rho \left(Qx^k - z^k \right) \right]^\top Qx \\
&\qquad\qquad\qquad (14) \\
&\geq \min_x \left\{ f(x) + \left(\omega + \rho \left(Qx^k - z^k \right) \right)^\top Qx : x \in \text{conv}(X) \right\}. \qquad (15)
\end{aligned}$$

Note that the term $-\rho(Qx^k - z^k)^\top z$ vanishes due to the optimality condition associated with (12). Inequality (13) follows from the inequalities (14)–(15) once the substitution $(x, z) = (x^k, z^k)$ and the definition of $\hat{\phi}(\omega, x^k, z^k)$ are applied to the left-hand side of (14). \square

Observe that, due to the linearity of the objective function with respect to x in (8), the use of constraint sets X and $\text{conv}(X)$ are interchangeable, and so in evaluating $\check{\phi}$, an explicit description of $\text{conv}(X)$ is not required. Furthermore, from the definition of ϕ^C , the convexity of f over \mathbb{R}^n , and the interchangeability of X and $\text{conv}(X)$ in (8), it is clear that for all $x^k \in \mathbb{R}^n$, $k \geq 0$, we have $\phi^C(\omega) \geq \check{\phi}(\omega, x^k)$. This is not only true in principle but also practically if one can provide an oracle that returns an extremal point of $\text{conv}(X)$ when minimizing a linear function over X . When the non-convexity is entirely due to the presence of integer restrictions on variables, MIP or MINLP solvers can provide such an oracle. Later in sect. 4 we shall see that the class of problems that is amenable to the final implementable algorithm is dictated, in practice, by the user's ability to provide such an oracle. Moreover, when f is linear, we have $\phi^C(\omega) \equiv \check{\phi}(\omega, x^k)$ for all x^k , $k \geq 0$; the two functions collapse into the same function with the centering at x^k of the latter function now irrelevant.

The first important property of $(\omega, x) \mapsto \check{\phi}(\omega, x)$ is its continuity.

Lemma 3 *Let X be compact, and f be continuously differentiable. Then $(\bar{\omega}, \bar{x}) \mapsto \check{\phi}(\bar{\omega}, \bar{x})$ is continuous over $(\bar{\omega}, \bar{x}) \in Z^\perp \times \mathbb{R}^n$.*

Proof From (8), compute

$$\begin{aligned}
\check{\phi}(\bar{\omega}, \bar{x}) &= f(\bar{x}) - \nabla_x f(\bar{x})^\top \bar{x} + \min_x \left\{ \left[\nabla_x f(\bar{x}) + \bar{\omega}^\top Q \right] x + \delta_{\text{conv}(X)}(x) \right\} \\
&= f(\bar{x}) - \nabla_x f(\bar{x})^\top \bar{x} - \delta_{\text{conv}(X)}^* \left(- \left[\nabla_x f(\bar{x}) + \bar{\omega}^\top Q \right] \right).
\end{aligned}$$

where $\delta_{\text{conv}(X)}(x) := \begin{cases} 0, & \text{if } x \in \text{conv}(X); \\ \infty, & \text{otherwise.} \end{cases}$ is the indicator function on the set $\text{conv}(X)$ and $\delta_{\text{conv}(X)}^*$ is the conjugate function [47] of $\delta_{\text{conv}(X)}$. As $\text{conv}(X)$ is convex and compact, we see that $\delta_{\text{conv}(X)}^*(\cdot)$ has domain \mathbb{R}^n and is thus continuous over \mathbb{R}^n (e.g., Lemma 2.91 of [4]), yielding the intended conclusion. \square

The second property of $\check{\phi}$ is its limiting behavior as the solutions (x^k, z^k) approach certain critical values.

Lemma 4 *Let the sequence $\{(x^k, z^k)\} \subset \text{conv}(X) \times Z$ satisfy the z -optimality condition (12) for each $k \geq 1$. If, for some fixed $\omega \in Z^\perp$, the sequence $\{(x^k, z^k)\}$ converges optimally in the sense that*

$$\lim_{k \rightarrow \infty} (x^k, z^k) = (x^*, z^*) \in \text{argmin}_{x, z} \{L_\rho(x, z, \omega) : x \in \text{conv}(X), z \in Z\},$$

then

$$\lim_{k \rightarrow \infty} \check{\phi}(\omega + \rho(Qx^k - z^k), x^k) = L_\rho(x^*, z^*, \omega) + \frac{\rho}{2} \|Qx^* - z^*\|_2^2. \quad (16)$$

Proof We begin by writing the necessary (and sufficient) conditions associated with the optimality $(x^*, z^*) \in \text{argmin}_{x, z} \{L_\rho(x, z, \omega) : x \in \text{conv}(X), z \in Z\}$:

$$\begin{bmatrix} \nabla f(x^*) + [\omega + \rho(Qx^* - z^*)]^\top Q \\ -\rho(Qx^* - z^*) \end{bmatrix} \begin{bmatrix} x - x^* \\ z - z^* \end{bmatrix} \geq 0 \quad \text{for all } x \in \text{conv}(X), z \in Z.$$

Since $z^k \in \text{argmin}_z \{\|Qx^k - z\|_2^2 : z \in Z\}$ for each $k \geq 1$, we have $Qx^k - z^k \in Z^\perp$, and so $Qx^* - z^* \in Z^\perp$ also. Thus, we can simplify the consideration of the above displayed necessary conditions to consider the x block only:

$$[\nabla f(x^*) + [\omega + \rho(Qx^* - z^*)]^\top Q]^\top [x - x^*] \geq 0 \quad \text{for all } x \in \text{conv}(X),$$

which implies

$$\min_x \left\{ [\nabla f(x^*) + [\omega + \rho(Qx^* - z^*)]^\top Q]^\top [x - x^*] : x \in \text{conv}(X) \right\} = 0.$$

In terms of $\check{\phi}(\omega + \rho(Qx^* - z^*), x^*)$, the above equality is re-written as:

$$\begin{aligned} \check{\phi}(\omega + \rho(Qx^* - z^*), x^*) &= f(x^*) + \omega^\top Qx^* + \rho \|Qx^* - z^*\|_2^2 \\ &= L_\rho(x^*, z^*, \omega) + \frac{\rho}{2} \|Qx^* - z^*\|_2^2, \end{aligned}$$

where the equality $(Qx^* - z^*)^\top z^* = 0$ is utilized. The continuity of $(\bar{\omega}, \bar{x}) \mapsto \check{\phi}(\bar{\omega}, \bar{x})$ established in Lemma 3 gives the desired conclusion. \square

We use Lemmas 3 and 4 to develop the proximal bundle method-like serious step condition (SSC) that makes use of $\hat{\phi}$ and $\check{\phi}$ in place of the cutting plane model and ϕ , respectively. Defining $\tilde{\omega}^k := \omega^k + \rho(Qx^k - z^k)$, consider the following modified serious step condition:

$$\gamma \leq \frac{\check{\phi}(\tilde{\omega}^k, x^k) - \check{\phi}(\omega^k, x^{k-1})}{\hat{\phi}(\omega^k, x^k, z^k) - \check{\phi}(\omega^k, x^{k-1})} \leq 1, \quad (17)$$

where $\gamma \in (0, 1)$ is the SSC parameter. The upper bound of (17) is satisfied automatically since $\widehat{\phi}(\omega^k, x^k, z^k) \geq \phi^C(\widetilde{\omega}^k) \geq \check{\phi}(\widetilde{\omega}^k, x^k)$ holds by Lemma 2 and the definition of $\check{\phi}$. However, the satisfaction of the lower bound is conditional on γ .

Remark 2 Throughout this paper, we shall always assume or construct z^k such that the z -optimality condition (12) is satisfied for each $k \geq 0$. Due to the necessary conditions of optimality associated with (12) and that Z is a linear subspace, we have $(Qx^k - z^k)^\top z = 0$ for all $z \in Z$. It immediately follows that if $\omega^k \in Z^\perp$, then $\widetilde{\omega}^k = \omega^k + \rho(Qx^k - z^k) \in Z^\perp$ also. Thus, the satisfaction of the z -optimality condition (12) guides the generation of $\{\omega^k\}$ so that if $\omega^0 \in Z^\perp$, then $\omega^k \in Z^\perp$ is always maintained for each $k \geq 1$.

Under certain circumstances, the denominator of the ratio displayed in (17) can be zero. The following lemma states that this never happens when ω^k is not dual optimal with respect to the dual problem (3).

Lemma 5 For any $\omega \in Z^\perp$ that is not dual optimal with respect to the dual problem (3) and $(x, z) \in \text{conv}(X) \times Z$, we have

$$\widehat{\phi}(\omega, x, z) - \phi^C(\omega) > 0. \quad (18)$$

Consequently, at any iteration k , the denominator of the ratio displayed in (17) cannot be zero when ω^k is not dual optimal.

Proof By the definition of $\widehat{\phi}$, we have

$$\widehat{\phi}(\omega, x, z) - \phi^C(\omega) \geq L_\rho(x^*, z^*, \omega) + \frac{\rho}{2} \|Qx - z\|_2^2 - \phi^C(\omega)$$

for all $\text{conv}(X)$ and $z \in Z$, where

$$(x^*, z^*) \in \text{argmin}_{x, z} \{L_\rho(x, z, \omega) : x \in \text{conv}(X), z \in Z\}.$$

(That is, we substitute $L_\rho(x, z, \omega)$ from the definition of $\widehat{\phi}$ with $L_\rho(x^*, z^*, \omega)$ to get the inequality.) Now $L_\rho(x^*, z^*, \omega) - \phi^C(\omega) > 0$ when ω is not dual optimal. Otherwise, if $L_\rho(x^*, z^*, \omega) = \phi^C(\omega)$, then $Qx^* = z^*$ must hold, and (x^*, z^*, ω) is a Lagrangian saddle point for problem (2) with respect to the Lagrangian relaxation of the constraint $Qx = z$. This contradicts the non-dual optimality of ω . Thus, the strict inequality (18) is established.

In the context of (17) at iteration k , noting that $\phi^C(\omega^k) \geq \check{\phi}(\omega^k, x^{k-1})$, we substitute $(x, z) = (x^k, z^k)$ and $\omega = \omega^k$ in the strict inequality (18) and so the denominator in (17) is positive when ω^k is not dual optimal. \square

From Lemma 4, we have the following result regarding the satisfaction of condition (17).

Proposition 1 Let the sequence $\{(x^k, z^k)\} \subset \text{conv}(X) \times Z$ satisfy

$$z^k \in \text{argmin}_z \left\{ \|Qx^k - z\|_2^2 : z \in Z \right\}$$

for each $k \geq 1$. Furthermore, let $\omega \in Z^\perp$ and $\omega \notin \operatorname{argmax}_\omega \phi(\omega)$. If the sequence $\{(x^k, z^k)\}$ converges optimally in the sense that

$$\lim_{k \rightarrow \infty} (x^k, z^k) = (x^*, z^*) \in \operatorname{argmin}_{x, z} \{L_\rho(x, z, \omega) : x \in \operatorname{conv}(X), z \in Z\},$$

then condition (17) must be satisfied after a finite number of iterations.

Proof For all $(x^k, z^k) \in \operatorname{conv}(X) \times Z$ with $z^k \in \operatorname{argmin}_z \{\|Qx^k - z\|_2^2\}$, we have

$$\begin{aligned} \widehat{\phi}(\omega, x^k, z^k) &= L_\rho(x^k, z^k, \omega) + \frac{\rho}{2} \|Qx^k - z^k\|_2^2 \\ &\geq \phi^C(\omega + \rho(Qx^k - z^k)) \geq \check{\phi}(\omega + \rho(Qx^k - z^k), x^k), \end{aligned}$$

where the first inequality follows from the definition of $\widehat{\phi}$ and Lemma 2, and the second inequality follows readily from the definition of $\check{\phi}$, the subgradient inequality and the interchangeability of X and $\operatorname{conv}(X)$. By the assumption that ω is not dual optimal, the denominator of (17) cannot be zero by Lemma 5. It follows from the convergence in (16) implied by Lemma 4 that the ratio in (17) must approach 1, and so condition (17) must be satisfied after a finite number of iterations. \square

Consequently, unless the current ω^k is already dual optimal, there cannot be an infinite number of null-steps when using condition (17). Recall that we use k to count only serious steps.

Proposition 2 Assume that problem (3) has an optimal dual solution ω^* , and that for each $k \geq 1$, $\phi^C(\omega^k) < \phi^C(\omega^*)$. Also, assume that ρ and γ may vary with each iteration, defined by sequences $\{\rho_k\}$ and $\{\gamma_k\}$ such that $\rho_k > 0$ and $\gamma_k \in (0, 1)$, bounded strictly away from zero for all $k \geq 1$, and $\rho_k \left(\frac{1-\gamma_k}{\gamma_k}\right) = c > 0$ for all k . If the sequence $\{\omega^k\}$ of dual updates is generated with Algorithm 1 with $\epsilon = 0$ and $k_{\max} = \infty$, then $\{\omega^k\}$ converges, and $\lim_{k \rightarrow \infty} \check{\phi}(\omega^k, x^{k-1}) = \zeta^{CLD}$ (and consequently $\lim_{k \rightarrow \infty} \phi^C(\omega^k) = \zeta^{CLD}$). Furthermore,

$$\lim_{k \rightarrow \infty} \widehat{\phi}(\omega^k, x^k, z^k) = \zeta^{CLD},$$

and all limit points (\bar{x}, \bar{z}) of the sequence $\{(x^k, z^k)\}$ are optimal for problem (2).

Proof Let ω^* be any dual optimal solution for problem (3). For each iteration $k \geq 1$, write the following two relations:

$$\begin{aligned} \|\omega^{k+1} - \omega^*\|_2^2 &= \|\omega^k - \omega^* + \rho_k(Qx^k - z^k)\|_2^2 \\ &= \|\omega^k - \omega^*\|_2^2 + 2\rho_k(Qx^k - z^k)^\top (\omega^k - \omega^*) + \rho_k^2 \|Qx^k - z^k\|_2^2, \end{aligned} \quad (19)$$

$$\begin{aligned} \text{and } \phi^C(\omega^*) &\leq L_{\rho_k}(x^k, z^k, \omega^*) = L_{\rho_k}(x^k, z^k, \omega^k) + (\omega^* - \omega^k)^\top (Qx^k - z^k) \\ &\implies (\omega^k - \omega^*)^\top (Qx^k - z^k) \leq L_{\rho_k}(x^k, z^k, \omega^k) - \phi^C(\omega^*). \end{aligned} \quad (20)$$

Substituting the inequality (20) into equality (19), we have

$$\begin{aligned} \|\omega^{k+1} - \omega^*\|_2^2 &\leq \|\omega^k - \omega^*\|_2^2 \\ &\quad + 2\rho_k \left[L_{\rho_k}(x^k, z^k, \omega^k) - \phi^C(\omega^*) \right] + \rho_k^2 \|Qx^k - z^k\|_2^2 \end{aligned} \quad (21)$$

$$\begin{aligned} &= \|\omega^k - \omega^*\|_2^2 + 2\rho_k \left[\check{\phi}(\omega^k, x^{k-1}) - \phi^C(\omega^*) \right] \\ &\quad + 2\rho_k \left[L_{\rho_k}(x^k, z^k, \omega^k) + \frac{\rho_k}{2} \|Qx^k - z^k\|_2^2 - \check{\phi}(\omega^k, x^{k-1}) \right]. \end{aligned} \quad (22)$$

By assumption, for each $k \geq 1$, we have $\phi^C(\omega^k) < \phi^C(\omega^*)$, so by Lemma 5 and $\epsilon = 0$, the Line 8 condition of Algorithm 1 never holds. Thus, the Line 9 condition, which is equivalent to the satisfaction of condition (17), is satisfied for each $k \geq 1$. Rewriting (17), with the substitution $\tilde{\omega}^k = \omega^{k+1}$, as

$$L_{\rho_k}(x^k, z^k, \omega^k) + \frac{\rho_k}{2} \|Qx^k - z^k\|_2^2 - \check{\phi}(\omega^k, x^{k-1}) \leq \frac{\check{\phi}(\omega^{k+1}, x^k) - \check{\phi}(\omega^k, x^{k-1})}{\gamma_k} \quad (23)$$

and substituting (23) into (22), we have

$$\begin{aligned} \|\omega^{k+1} - \omega^*\|_2^2 &\leq \|\omega^k - \omega^*\|_2^2 + 2\rho_k \left[\check{\phi}(\omega^k, x^{k-1}) - \phi^C(\omega^*) \right] \\ &\quad + \frac{2\rho_k}{\gamma_k} \left[\check{\phi}(\omega^{k+1}, x^k) - \check{\phi}(\omega^k, x^{k-1}) \right] \\ &\leq \|\omega^k - \omega^*\|_2^2 + 2\rho_k \left[\check{\phi}(\omega^{k+1}, x^k) - \phi^C(\omega^*) \right] \\ &\quad + 2\rho_k \left(\frac{1 - \gamma_k}{\gamma_k} \right) \left[\check{\phi}(\omega^{k+1}, x^k) - \check{\phi}(\omega^k, x^{k-1}) \right] \end{aligned} \quad (24)$$

From (24), we make the following three inferences: (1) that $\{\|\omega^k - \omega^*\|\}$ is bounded, (2) that $\sum_{k=1}^{\infty} [\phi^C(\omega^*) - \phi^C(\omega^k)]$ is finite, and (3) that $\{\omega^k\}$ converges. To establish these inferences, we rearrange terms and sum the inequality (24) from $k = \ell, \dots, N$ for some integers $1 \leq \ell \leq N$ to get

$$\begin{aligned} &2 \sum_{k=\ell}^N \rho_k \left[\phi^C(\omega^*) - \check{\phi}(\omega^{k+1}, x^k) \right] + \|\omega^{N+1} - \omega^*\|_2^2 \\ &\leq \|\omega^\ell - \omega^*\|_2^2 + 2\rho_N \left(\frac{1 - \gamma_N}{\gamma_N} \right) \check{\phi}(\omega^{N+1}, x^N) - 2\rho_\ell \left(\frac{1 - \gamma_\ell}{\gamma_\ell} \right) \check{\phi}(\omega^\ell, x^{\ell-1}) \end{aligned}$$

$$\begin{aligned}
& + 2 \sum_{k=\ell}^{N-1} \left[\rho_k \left(\frac{1-\gamma_k}{\gamma_k} \right) - \rho_{k+1} \left(\frac{1-\gamma_{k+1}}{\gamma_{k+1}} \right) \right] \check{\phi}(\omega^{k+1}, x^k) \\
& \leq \|\omega^\ell - \omega^*\|_2^2 + 2c \left[\phi^C(\omega^*) - \check{\phi}(\omega^\ell, x^{\ell-1}) \right], \tag{25}
\end{aligned}$$

where the last inequality follows from the assumption that $\rho_k \left(\frac{1-\gamma_k}{\gamma_k} \right) = c$ for all k and the bounding relationships implied by the optimality of ω^* :

$$\phi^C(\omega^*) > \phi^C(\omega^{k+1}) \geq \check{\phi}(\omega^{k+1}, x^k) \tag{26}$$

Noting that each summand $\phi^C(\omega^*) - \check{\phi}(\omega^{k+1}, x^k)$ in the summation on the left-hand side of (25) is nonnegative so we have immediately from (25) that $\sum_{k=1}^{\infty} \rho_k \left[\phi^C(\omega^*) - \check{\phi}(\omega^k, x^{k-1}) \right] < \infty$ and $\{\omega^k - \omega^*\}$ is bounded, establishing the first two inferences from (24). The validity of the first two inferences imply the boundedness of $\{\omega^k\}$ and the convergence $\lim_{k \rightarrow \infty} \check{\phi}(\omega^k, x^{k-1}) = \phi^C(\omega^*)$, respectively. The boundedness of $\{\omega^k\}$ implies the existence of limit points, while the convergence $\lim_{k \rightarrow \infty} \check{\phi}(\omega^k, x^{k-1}) = \phi^C(\omega^*)$ implies that all such limit points are dual optimal. It is straightforward from (26) that $\lim_{k \rightarrow \infty} \phi^C(\omega^k) = \phi^C(\omega^*)$ also.

To establish the third assertion, that $\{\omega^k\}$ in fact converges, we drop the summation from the left-hand side of (25),

$$\|\omega^{N+1} - \omega^*\|_2^2 \leq \|\omega^\ell - \omega^*\|_2^2 + 2c \left[\phi^C(\omega^*) - \check{\phi}(\omega^\ell, x^{\ell-1}) \right], \tag{27}$$

and note that the above analysis holds independent of the choice of dual optimal ω^* . Since it was just shown that $\{\omega^k\}$ has limit points, and that all such limit points are dual optimal, we now specify ω^* to be one of these limit points. We then choose an appropriate ℓ for any $\varepsilon > 0$ so that the right-hand side of (27) is arbitrarily small, i.e.,

$$\|\omega^{N+1} - \omega^*\|_2^2 \leq \varepsilon$$

for all $N \geq \ell$. Thus, $\lim_{k \rightarrow \infty} \omega^k = \omega^*$, and it is clear that the limit point ω^* of $\{\omega^k\}$ is in fact unique.

To prove the last assertion, the satisfaction of (17) is rewritten as

$$\begin{aligned}
\check{\phi}(\omega^{k+1}, x^k) - \check{\phi}(\omega^k, x^{k-1}) & \leq \widehat{\phi}(\omega^k, x^k, z^k) - \check{\phi}(\omega^k, x^{k-1}) \\
& \leq \frac{1}{\gamma_k} \left(\check{\phi}(\omega^{k+1}, x^k) - \check{\phi}(\omega^k, x^{k-1}) \right).
\end{aligned}$$

Due to the convergence $\lim_{k \rightarrow \infty} \check{\phi}(\omega^k, x^{k-1}) = \zeta^{CLD}$, we have on taking the limit as $k \rightarrow \infty$ of the last displayed inequalities that $\lim_{k \rightarrow \infty} \widehat{\phi}(\omega^k, x^k, z^k) = \zeta^{CLD}$. Noting that $\widehat{\phi}(\omega^k, x^k, z^k) = L_{\rho_k}(x^k, z^k, \omega^k) + \frac{\rho_k}{2} \|Qx^k - z^k\|_2^2$, it is clear

that if $\limsup_{k \rightarrow \infty} \rho_k = \infty$, then we have $\lim_{k \rightarrow \infty} \|Qx^k - z^k\|_2^2 = 0$ and $\lim_{k \rightarrow \infty} L_{\rho_k}(x^k, z^k, \omega^k) = \zeta^{CLD}$, and so the limit points of $\{(x^k, z^k)\}$ must be feasible and furthermore optimal for (2). Now assume $0 < \limsup \rho_k < \infty$. In taking the limit points $(\bar{x}, \bar{z}, \omega^*)$ of the sequence $\{(x^k, z^k, \omega^k)\}$ and $\bar{\rho}$ of $\{\rho_k\}$, noting that the optimal value of problem (9) with $\omega = \omega^*$ is ζ^{CLD} by Lemma 1,

$$\zeta^{CLD} + \frac{\bar{\rho}}{2} \|Q\bar{x} - \bar{z}\|_2^2 \leq L_{\bar{\rho}}(\bar{x}, \bar{z}, \omega^*) + \frac{\bar{\rho}}{2} \|Q\bar{x} - \bar{z}\|_2^2 = \zeta^{CLD}.$$

From this, it follows that $\|Q\bar{x} - \bar{z}\|_2^2 = 0$ and $L_{\bar{\rho}}(\bar{x}, \bar{z}, \omega^*) = \zeta^{CLD}$, and so (\bar{x}, \bar{z}) must be feasible and furthermore optimal for (2). \square

3.1 Rate-of-convergence

The proof of Proposition 2 allows for the following remarks on the rate-of-convergence associated with $\lim_{k \rightarrow \infty} \check{\phi}(\omega^k, x^{k-1}) = \zeta^{CLD}$. Note that each iteration k of Algorithm 1 corresponds to a serious step update of ω^k .

1. Let $\rho_k = \rho$ and $\gamma_k = \gamma$ for all $k \geq 1$, where $\rho > 0$ and $\gamma \in (0, 1)$ are constants. Then, from the proof of Proposition 2, we have

$$\sum_{k=1}^N \left[\phi^C(\omega^*) - \check{\phi}(\omega^k, x^{k-1}) \right] < \infty$$

and since $\{\check{\phi}(\omega^k, x^{k-1})\}$ is monotonically non-decreasing, it is clear that

$$\phi^C(\omega^*) - \check{\phi}(\omega^k, x^{k-1}) = o(1/k),$$

where o is the Little-o notation.

2. Let $\rho_k = k\rho$ for some constant $\rho > 0$ and $\gamma_k = \frac{\rho_k}{c+\rho_k}$ for some constant $c > 0$ so that $\rho_k \left(\frac{1-\gamma_k}{\gamma_k} \right) = c$ is a constant. Then we have

$$\sum_{k=1}^N k \left[\phi^C(\omega^*) - \check{\phi}(\omega^k, x^{k-1}) \right] < \infty,$$

and so $\phi^C(\omega^*) - \check{\phi}(\omega^k, x^{k-1}) = o(1/k^2).$

3. Let $\rho_k = \rho b^k$ for some constant $\rho > 0$ and $b > 1$, and as $\gamma_k = \frac{\rho_k}{c+\rho_k}$ for some constant $c > 0$ so that $\rho_k \left(\frac{1-\gamma_k}{\gamma_k} \right) = c$, a constant. Then we have

$$\sum_{k=1}^N b^k \left[\phi^C(\omega^*) - \check{\phi}(\omega^k, x^{k-1}) \right] < \infty,$$

and so $\phi^C(\omega^*) - \check{\phi}(\omega^k, x^{k-1}) = o(b^{-k}).$

Since the number of null step updates per serious step is not fixed, and a null step does not require significantly less computational effort, these convergence results in terms of serious steps cannot be generalised to a convergence result in terms of total steps or runtime without some notion of how often null steps are taken.

Exploratory numerical tests did not conclusively reveal any clear and consistent pattern to the frequency of null steps as the algorithm progresses. However, for at least some combinations of instances and parameters the null step frequency increased significantly as the duality gap decreased. Therefore, the practical convergence behaviour of the algorithm in terms of runtime is likely worse than the above results would suggest.

4 Main algorithm

After integrating SDM and the nonlinear block Gauss–Seidel method, a practical implementation of Algorithm 1 is provided in this section.

We consider the following general two-block problem

$$\min_{x,z} \{F(x, z) : x \in \text{conv}(X), z \in Z\} \quad (28)$$

where $F : \mathbb{R}^n \times \mathbb{R}^q \mapsto \mathbb{R}$ is a continuously differentiable function, $\text{conv}(X)$ and Z are closed convex sets, and $\text{conv}(X)$ is also bounded. (Z can be more generally a convex set in this setting, not necessarily a linear (sub)space.) Additionally, we assume for each fixed $x \in \text{conv}(X)$ that $z \mapsto F(x, z)$ is inf-compact. (That is, the set $\{z \in Z : F(x, z) \leq \ell\}$ is compact for all $x \in \text{conv}(X)$ and $\ell \in \mathbb{R}$.) In the context of Algorithm 1, we would identify $F(x, z) = L_\rho(x, z, \omega)$ for a given ω .

Problem (28) is assumed to be feasible, bounded, and to have an optimal solution (x^*, z^*) . We shall utilize the following two-block nonlinear Gauss–Seidel (GS) method with the x update approximated in a manner resembling an iteration of the SDM. We assume the user provides an oracle to return an extremal point in $\text{conv}(X)$ when minimizing a linear function over X . This can be used to initialize the following and later algorithms.

Algorithm 2 An iteration of inner-approximated nonlinear Gauss–Seidel approach applied to problem (28).

```

1: Precondition:  $\tilde{x} \in \text{conv}(X)$ ,  $\tilde{z} \in \text{argmin}_z \{F(\tilde{x}, z) : z \in Z\}$ ,  $D \subseteq \text{conv}(X)$ 
2: function SDM- GS( $F, X, Z, D, \tilde{x}, \tilde{z}, t_{\max}$ )
3:   for  $t = 1, \dots, t_{\max}$  do
4:      $\tilde{x} \leftarrow \text{argmin}_x \{F(x, \tilde{z}) : x \in D\}$ 
5:      $\tilde{z} \leftarrow \text{argmin}_z \{F(\tilde{x}, z) : z \in Z\}$ 
6:   end for
7:    $\hat{x} \in \text{argmin}_x \left\{ \nabla_x F(\tilde{x}, \tilde{z})^\top (x - \tilde{x}) : x \in X \right\}$ 
8:   Reconstruct  $D$  to be any set such that
9:      $\{\tilde{x} + \alpha(\hat{x} - \tilde{x}) : \alpha \in [0, 1]\} \subseteq D \subseteq \text{conv}(X)$ 
10:   Set  $\Gamma \leftarrow -\nabla_x F(\tilde{x}, \tilde{z})(\hat{x} - \tilde{x})$ 
11:   return  $(\tilde{x}, \tilde{z}, D, \Gamma)$ 
12: end function

```

If the z block update of Line 5 is trivialized, such as by making z not actually appear in the definition of F , or by making Z a singleton set, then Algorithm 2 would be identical to SDM applied to problem (28) in which the z block of variables correspondingly does not play any role. On the other hand, if the x update (4) is replaced with an update based on an exact minimization $\tilde{x} \leftarrow \operatorname{argmin}_x \{F(x, \tilde{z}) : x \in \operatorname{conv}(X)\}$ (so that the computations of Lines 7–10 and the returning of D and Γ can be skipped), then Algorithm 2 would be equivalent to a more traditional two-block nonlinear Gauss–Seidel method. Different forms of approximation of the x update, such as those resulting from gradient descent steps in x , are also considered in [21, 48].

Remark 3 The main approach envisioned for constructing the inner approximation D on Lines 8–9 is to take $D \leftarrow \operatorname{conv}(D \cup \{\tilde{x}, \hat{x}\})$. To implement this update of D , we need to save the points \hat{x} computed during previous calls to Algorithm 2.

We assume in the following proposition that Algorithm 2 is applied iteratively in the sense that at iteration $k \geq 0$, we input $(\tilde{x}, \tilde{z}) = (x^k, z^k)$ and return $(\tilde{x}, \tilde{z}) = (x^{k+1}, z^{k+1})$. Furthermore, at the same iteration k call of Algorithm 2, we set $d^{k+1} = \hat{x} - \tilde{x}$ where \hat{x} and \tilde{x} are set as in Line 9. This provides a reference sequence of directions $\{d^k\}$ necessary in the proof of the following proposition.

Proposition 3 *For problem (28), let F be convex and continuously differentiable, and let $\operatorname{conv}(X)$ and Z be nonempty and convex, with $\operatorname{conv}(X)$ bounded and $z \mapsto F(x, z)$ inf-compact for each $x \in \operatorname{conv}(X)$. Then, for any $t_{\max} \geq 1$, the sequence $\{(x^k, z^k)\}$ generated by iterations of Algorithm 2 has limit points (\bar{x}, \bar{z}) , each of which are optimal for problem (28).*

Proof In light of the convexity and continuous differentiability of F and the convexity of $\operatorname{conv}(X)$ and Z , it is sufficient to show that

$$\nabla_x F(\bar{x}, \bar{z})^\top (x - \bar{x}) \geq 0 \quad \text{for all } x \in \operatorname{conv}(X) \quad (29)$$

$$\text{and } \nabla_z F(\bar{x}, \bar{z})^\top (z - \bar{z}) \geq 0 \quad \text{for all } z \in Z. \quad (30)$$

As $\nabla_z F(x^k, z^k)^\top (z - z^k) \geq 0$ for all $z \in Z$ holds for each $k \geq 1$ (this follows due to the optimality $z^k \in \operatorname{argmin}_z \{F(x^k, z) : z \in Z\}$ that holds by construction) the satisfaction of the latter condition (30) is trivially established for any limit points (\bar{x}, \bar{z}) . It remains only to show the satisfaction of the x -stationarity condition (29). This may be established by using Proposition 3.2 of [21] combined with the last sentence of Remark 3.3 from the same reference. But for the sake of explicitness, we use developments in Appendix A to show that (29) holds.

Note, for the sake of nontriviality, that $\nabla_x F(x^k, z^k)^\top (x - x^k) \geq 0$ for all $x \in X$ is assumed *not* to hold for any $k \geq 1$. Thus, with reference to the argument given in Appendix A, the sequence of directions $\{d^k\}$ satisfy the the Direction Assumption (DA), prior to Algorithm 4. Also, the Gradient Related Assumption (GRA) referred to in Appendix A is satisfied for this same $\{d^k\}$, by Lemma 7 therein. Due to the construction of D in Line 9 and setting $(x^{k+1}, z^{k+1}) = (\tilde{x}, \tilde{z})$ after the termination of the for loop of Lines 3–6, we have given $\{d^k\}$ and any choice of $(\beta, \sigma) \in (0, 1)$ the satisfaction of the Sufficient Decrease Assumption (SDA), also referred to in

Appendix A. It then follows from Lemma 6 that limit points (\bar{x}, \bar{z}) of $\{(x^k, z^k)\}$ do exist, and that each of which satisfies the stationarity condition (29). \square

The method SDM-GS-ALM is now stated as Algorithm 3, which uses Algorithm 2 as a subroutine to provide a practical implementation of Algorithm 1.

Remark 4 At the return of Algorithm 2 in Line 8 of Algorithm 3, we have

$$\begin{aligned}\Gamma &= -\nabla_x L_\rho(x^k, z^k, \omega^k)^\top (\hat{x} - x^k) \\ &= -\left[\nabla_x f(x^k) + \left(\omega^k + \rho(Qx^k - z^k) \right)^\top Q \right]^\top (\hat{x} - x^k)\end{aligned}$$

where \hat{x} is computed on Line 7 of Algorithm 2. One may verify the equality $(Qx^k - z^k)^\top z^k = 0$ due to $z^k \in \operatorname{argmin}_z \{\|Qx^k - z\|_2^2 : z \in Z\}$. Moreover using this value of Γ and the computation of $\tilde{\phi}$ on Lines 4 and 12 one may show, using the fact that $\hat{x} \in \operatorname{argmin}_x \{\nabla_x L_\rho(x^k, z^k, \omega^k)^\top (x - x^k) : x \in X\}$, that

$$\tilde{\phi} = L_\rho(x^k, z^k, \omega^k) + \frac{\rho}{2} \|Qx^k - z^k\|_2^2 - \Gamma = \check{\phi}(\omega^k + \rho(Qx^k - z^k), x^k).$$

Algorithm 3 A practical implementation of Algorithm 1 based on the use of SDM-GS iterations. (SDM-GS is given as Algorithm 2.)

```

1: Preconditions:  $x^0 \in \operatorname{conv}(X)$ ,  $z^0 \in Z$ ,  $\omega^0 \in Z^\perp$ ,  $D \subseteq \operatorname{conv}(X)$ ,  $\gamma \in (0, 1)$ .
2: function SDM-GS-ALM( $f, Q, X, Z, D, \rho, x^0, z^0, \omega^0, \gamma, \epsilon, t_{\max}, k_{\max}$ )
3:    $(x^0, z^0, D, \Gamma) \leftarrow \text{SDM-GS}(L_\rho(\cdot, \cdot, \omega^0), X, Z, D, x^0, z^0, t_{\max})$ 
4:    $\tilde{\phi} \leftarrow L_\rho(x^0, z^0, \omega^0) + \frac{\rho}{2} \|Qx^0 - z^0\|_2^2 - \Gamma$ 
5:   set  $\omega^0 \leftarrow \omega^0 + \rho(Qx^0 - z^0)$ ,  $\check{\phi}^0 \leftarrow \tilde{\phi}$ 
6:   for  $k = 1, 2, \dots, k_{\max}$  do
7:     Initialize  $\omega^k \leftarrow \omega^{k-1}$ ,  $\check{\phi}^k \leftarrow \check{\phi}^{k-1}$  ▷ (Default, null-step updates)
8:      $(x^k, z^k, D, \Gamma) \leftarrow \text{SDM-GS}(L_\rho(\cdot, \cdot, \omega^k), X, Z, D, x^{k-1}, z^{k-1}, t_{\max})$ 
9:     if  $L_\rho(x^k, z^k, \omega^k) + \frac{\rho}{2} \|Qx^k - z^k\|_2^2 - \check{\phi}^k \leq \epsilon$  then
10:      return  $(x^k, z^k, \omega^k, \check{\phi}^k)$ 
11:    end if
12:     $\tilde{\phi} \leftarrow L_\rho(x^k, z^k, \omega^k) + \frac{\rho}{2} \|Qx^k - z^k\|_2^2 - \Gamma$ 
13:     $\eta_k \leftarrow \frac{\tilde{\phi} - \check{\phi}^k}{L_\rho(x^k, z^k, \omega^k) + \frac{\rho}{2} \|Qx^k - z^k\|_2^2 - \check{\phi}^k}$ 
14:    if  $\eta_k \geq \gamma$  then
15:      set  $\omega^k \leftarrow \omega^k + \rho(Qx^k - z^k)$ ,  $\check{\phi}^k \leftarrow \tilde{\phi}$ 
16:    end if
17:    Possibly update  $\rho$ , e.g.,  $\rho \leftarrow \frac{1}{\min\{\max\{(2/\rho)(1-\eta_k), 1/(10\rho), 10^{-4}\}, 10/\rho\}}$  as in [49]
18:  end for
19:  return  $(x^k, z^k, \omega^k, \check{\phi}^k)$ 
20: end function

```

Proposition 4 Let $\{(x^k, z^k, \omega^k)\}$ be a sequence generated by Algorithm 3 applied to problem (1) with X compact, Z a linear subspace, $\omega^0 \in Z^\perp$, $\rho > 0$, $\gamma \in (0, 1)$, $\epsilon = 0$ and $k_{\max} = \infty$. If there exists a dual optimal solution ω^* to the dual problem (3), then either

1. $\omega^k = \bar{\omega}$ is fixed and optimal for (3) for $k \geq \bar{k}$ for some finite \bar{k} ; or
2. ω^k is never optimal for (3) for any finite $k \geq 1$, but $\lim_{k \rightarrow \infty} \omega^k = \bar{\omega}$ is optimal,

and the sequence $\{(x^k, z^k)\}$ has limit points (\bar{x}, \bar{z}) , each of which is optimal for problem (2).

Proof In the first case, Algorithm 3 never takes serious steps for iterations $k \geq \bar{k} \geq 1$, and so with $\omega^k = \bar{\omega}$, optimal for (3) and fixed for $k \geq \bar{k}$, the Algorithm 3 iterations continue with the generation of $\{(x^k, z^k)\}$ as generated by iterations of SDM-GS (Algorithm 2). By Proposition 3, the sequence $\{(x^k, z^k)\}$ has limit points (\bar{x}, \bar{z}) , each of which is optimal for problem (9) with $\omega = \bar{\omega}$. Then, by Lemma 1, (\bar{x}, \bar{z}) is also optimal for problem (2) since $\bar{\omega}$ is optimal for (3).

In the second case, where ω^k is never dual optimal for (2) for any finite $k \geq 1$, any serious step must be followed by a finite number of consecutive null-steps. We consider the subsequence indices $\{k_i\}_{i=1}^\infty$ where the update ω^{k_i+1} is obtained by a serious step. By Proposition 2, we have $\lim_{i \rightarrow \infty} \phi^C(\omega^{k_i+1}) = \zeta^{CLD}$, and accommodating the null steps in between, we have also $\lim_k \phi^C(\omega^k) = \zeta^{CLD}$. To prove the last claim, we note that $\omega^j = \omega^{k_i+1}$ for all integers j such that $k_i < j \leq k_{i+1}$ due to the taking of null steps. From Proposition 2, we have that $\lim_{i \rightarrow \infty} L_\rho(x^{k_i}, z^{k_i}, \omega^{k_i}) = \zeta^{CLD}$. By the continuity of $(x, z, \omega) \mapsto L_\rho(x, z, \omega)$, the convergences $\lim_{k \rightarrow \infty} \omega^k = \bar{\omega}$ and $\lim_{i \rightarrow \infty} Qx^{k_i} - z^{k_i} = 0$ (again, Proposition 2), we have $\lim_{i \rightarrow \infty} L_\rho(x^{k_i}, z^{k_i}, \omega^{k_i+1}) = \zeta^{CLD}$ also. Next, at each i , and integers j such that $k_i < j \leq k_{i+1}$, observe that

$$L_\rho(x^{k_i}, z^{k_i}, \omega^{k_i+1}) \geq L_\rho(x^j, z^j, \omega^j) \geq L_\rho(x^{k_{i+1}}, z^{k_{i+1}}, \omega^{k_{i+1}}).$$

In taking the limit of the above inequality as $i \rightarrow \infty$, it becomes evident that $\lim_{k \rightarrow \infty} L_\rho(x^k, z^k, \omega^k) = \zeta^{CLD}$ in the original sequence also. By the optimality of $\bar{\omega}$ for problem (3), we know from Lemma 1 that $\zeta_\rho^{AL}(\bar{\omega}) = \zeta^{CLD}$, and so each limit point (\bar{x}, \bar{z}) must be optimal for problem (9) with $\omega = \bar{\omega}$. Furthermore, by Lemma 1, (\bar{x}, \bar{z}) must also be optimal for problem (2). (These limit points exist furthermore, due to the compactness of $\text{conv}(X)$ and the continuous and closed-form expression that the unique solution $z^k \in \arg\min_z \{\|Qx^k - z\|_2^2 : z \in Z\}$ has given $x^k \in \text{conv}(X)$ when Z is a linear subspace.) \square

4.1 Parallelization and workload

The opportunities for parallelization and distribution of the computational workload in SDM-GS-ALM, as stated in Algorithm 3, are not immediately apparent. This subsection explicitly indicates which update problems may be solved in parallel, and the nature of the required communication between the parallel computational nodes.

The bulk of computational work, parallelization, and parallel communication occurs within the SDM-GS method stated in Algorithm 2, where for the problems of interest, the following decomposable structures apply: $X = \prod_{i=1}^m X_i$, $D = \prod_{i=1}^m D_i$, and $F(x, z) = \sum_{i=1}^m F(x_i, z)$. In the larger context of Algorithm 3, the subproblem of Line 4 in Algorithm 2 can be solved in parallel given fixed $\tilde{z} \in Z$ and $\omega \in Z^\perp$ along the block indices $i = 1, \dots, m$ as

$$\min_x \left\{ f_i(x) + \omega_i^\top Q_i x + \frac{\rho}{2} \|Q_i x - \tilde{z}_i\|_2^2 : x \in D_i \right\}, \quad (31)$$

while the subproblem of Line 7 is solved as

$$\min_x \left\{ \nabla_x f_i(\tilde{x}_i) + (\omega_i + \rho(Q_i \tilde{x}_i - \tilde{z}_i))^\top Q_i x : x \in X_i \right\}.$$

Remark 5 In the setting where problem (1) is a large-scale mixed-integer linear optimization problem, the subproblems of Line 4 are continuous convex quadratic optimization problems for each block $i = 1, \dots, m$, which can be solved independently of one another and in parallel. In the same setting, the Line 7 subproblems are mixed-integer optimization problems for each block $i = 1, \dots, m$, which can also be solved independently of one another and in parallel. Additionally, the reconstruction of D occurring in Line 9 can be done in parallel for each D_i along the indices $i = 1, \dots, m$.

Parallel communication is needed for the computation of the z update in Line 5 in Algorithm 2. In the larger context of Algorithm 3, this takes the form of solving

$$\min_z \left\{ \sum_{i=1}^m \|Q_i \tilde{x}_i - z_i\|_2^2 : z \in Z \right\}.$$

This is solved as an averaging that requires the reduce-sum type parallel communication. The evaluation of the serious step condition through calculating η_k in Line 13 in Algorithm 3 also requires a reduce-sum type parallel communication. For implementation purposes, the computation of these values, including the computation of Γ from the SDM-GS call, can be combined into one reduce-sum communication. In total, each iteration of Algorithm 3 requires two reduce-sum type communications, one for computing the z -update of Line 5 Algorithm 2, and one combined reduce-sum communication to compute scalars associated with the Lagrangian bounds and the critical values for the termination conditions. The storage and updates of x^k and ω^k and D can also be done in parallel, while z^k and η_k need to be computed and stored by every processor at each iteration k .

5 Computational experiments and results

In this section, we present and examine the results of two computational tests with the following purposes:

- Test 1:** to demonstrate the effect of enforcing the serious step condition on the Lagrangian values;
- Test 2:** to compare the parallel speedup obtained with the use of two parallel implementations of SDM-GS-ALM (Algorithm 3) versus the parallel speedups reported in [17] for two other parallel approaches. Additionally, the final iteration Lagrangian bounds are compared between the different parallel implementations for each experiment.

Computational experiments were performed on instances from three classes of problems. The first class consists of the capacitated allocation problems (CAP) [50]. The second and third classes consist of the DCAP and SSLP problems from the Stochastic Integer Programming Test Problem Library (SIPLIB), which are described in detail in [51, 52] and accessible at [51]. These are all large-scale mixed-integer linear optimization problems, so the preceding observations for when f is linear apply.

Test 1 was conducted with a Matlab 2012b [53] serial implementation of Algorithm 3 using CPLEX 12.6.1 [54] as the solver. The computing environment was on an Intel® Core™ i7-4770 3.40 GHz processor with 8 GB RAM and on a 64-bit operating system. All experiments for Test 1 were run with maximum number of iterations $k_{max} = 20$.

The parallel experiments of Test 2 were conducted with a C++ implementation of Algorithm 3 using CPLEX 12.5 [55] as the solver and the message passing interface (MPI) for parallel communication. For reading SMPS files into scenario-specific subproblems and for their interface with CPLEX, we used modified versions of the COIN-OR [56] Smi and Osi libraries, either to instantiate appropriate C++ class instances of the subproblems directly, or to write scenario-specific MPS files from the SMPS file. The computing environment for the Test 2 experiments is the Raijin cluster maintained by Australia's National Computing Infrastructure (NCI) and supported by the Australian government [57]. The Raijin cluster is a high performance computing (HPC) environment which has 3592 nodes (system units), 57472 cores of Intel Xeon E5-2670 processors with up to 8 GB PC1600 memory per core (128 GB per node). All experiments were conducted using one thread per CPLEX solve.

5.1 Effects of the serious step condition

The results of the Test 1 set of experiments are depicted in the plots of Fig. 1 (with additional Figures 2, 3 and 4 in Appendix B). The use of different penalty parameter ρ values is differentiated by the use of different plot colors. The penalties are chosen so that the smallest penalties (in red) are near optimal in terms of the resulting computational performance, while the larger penalties are known beforehand to be too large for optimal performance. For testing purposes, this is the most interesting way to choose penalty values, as smaller (than optimal) penalty values yield very little difference in Lagrangian bound between the use of different SSC parameter values. Solid line and dashed line plots depict the Lagrange bounds due to the use of a more stringent SSC parameter value $\gamma = 0.5$ and a more lenient value for the SSC parameter $\gamma = 0.125$, respectively. The dotted line plots depict the Lagrangian values resulting

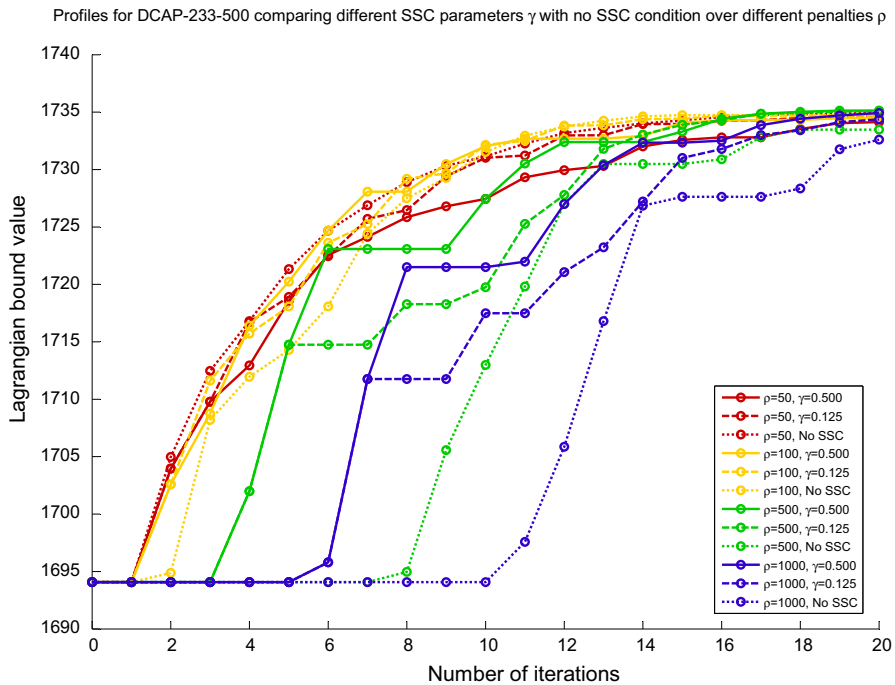


Fig. 1 Applying SDM-GS-ALM using different parameterizations for the SSC condition (or none)

from the non-use of the SSC, so that it evaluates true no matter what. The following observations are suggested from the results of the Test 1 experiments:

1. First, the most significant differences between the varied use of SSC occur when the penalty coefficient values are large. In this setting, it seems to be the case that the use of more stringent (i.e., larger) values of the SSC parameter γ has the effect of mitigating the destabilizing effect of having a penalty parameter ρ value that is too large. This is significant because the performance of iterative Lagrangian dual solution approaches based on (or related to) proximal bundle methods is sensitive to the tuning of the ρ value, and the optimal tuning of such parameters is assumed to be unknown beforehand in practical applications. For this reason, any mechanism to mitigate the effect of having an unfavorable tuning of the penalty parameter is highly desirable.
2. As is the case for the proximal bundle method, information from the SSC test can be used to dynamically fine-tune the value of the penalty parameter ρ .
3. While not enforcing the SSC can adversely affect the growth trend in the Lagrangian bound, the use of a SSC parameter γ value that is too large can have a similar effect for the tail-end values. This is most clearly seen in the Fig. 1 DCAP-233-500 $\rho = 50$ and $\rho = 100$ plots. In these plots, the growth in Lagrangian bound value is noticeably stunted in the tail-end iterations for the larger $\gamma = 0.5$ value as compared with the smaller $\gamma = 0.125$.

5.2 Benefits of parallelization

For the Test 2 experiments, we primarily compare the parallel speedup achieved with Algorithm 3 against that achieved with the enhancements to the proximal bundle method presented in [17]. Additionally, we compare the Lagrangian bound at the final iteration.

The enhancements in [17] use structure-exploiting primal-dual interior point solvers to improve the parallel efficiency of solving the proximal bundle method master problem. (The solution of this master problem is analogous to the approximated solution to problem (9) obtained by using the SDM-GS method in Algorithm 2.) The first solver is referred to by its acronym OOQP [58], while the second is PIPS-IPM [59].

In the experiments of Test 2, the underlying computing architecture and third-party software are inevitably different between our tests and those in [17]. Additionally, the termination criterion is necessarily different from that given in Step 2 of Figure 2 in [17] due to the differences in algorithms. In our tests, the termination criterion comes from Lines 9–11 of Algorithm 3 with $\epsilon = 10^{-6}$. We can nevertheless create a meaningful control in the tuning of the most important parameters affecting the performance of the algorithm.

1. As done in [17], we set the SSC parameter $\gamma = 0.1$, and we initialize the dual solution $\omega^0 = 0$.
2. In analogy to the possible trimming of cutting planes noted in [17], practical implementations of Algorithm 3 may judiciously trim the set D to improve performance. As all cuts are kept in the experiments of [17], so we also avoid trimming the expansion of D in our experiments, and so we just use the simple update rule $D \leftarrow \text{conv}(D \cup \{\tilde{x}, \hat{x}\})$ within Algorithm 2.
3. We use an update rule analogous to the one in [49] as is done in [17], which takes the suggested form given in Line 17 of Algorithm 3. Initially, $\rho = 1$.

In Tables 1 and 2, the columns headed by OOQP and PIPS-IPM report the parallel speedup due to the use of $N = 1, 8, 16, 32$ processors, which are originally reported in Figure 2 of [17]. If, given the use of N processors, T_N denotes the total wall clock time (in seconds) divided by number of iterations, then we compute the parallel speedup as T_1/T_N . For the computational experiments with Algorithm 3, we compute each table entry T_1/T_N after taking, from five identically parameterized experiments, (1) the minimum T_1 value, and (2) the average $T_N, N > 1$, value. The column headed by SDM-GS1-ALM presents the parallel speedup values for the application of Algorithm 3 with $t_{max} = 1$. The column headed by SDM-GS5-ALM is analogous, with $t_{max} = 5$. The total wall clock time per iteration values used to compute the ratios T_1/T_N are provided in Appendix C, accounting for taking the minimum ($N = 1$) or average ($N > 1$) over the five experiments for each set of parameterizations associated with Algorithm 3. For the two sets of experiments based on the application of Algorithm 3, a problem-specific maximum number of main loop iterations was set so as to make the tests as comparable with the tests in [17] as possible. These data are also reported in Appendix C. Also in Tables 1 and 2, the best Lagrangian bounds obtained for each combination of test problem and algorithm are reported.

Table 1 SSLP: comparing speedup and final best Lagrangian bound

No. Proc.	Speedup for SSLP 5-25-100			
	OOQP	PIPS-IPM	SDM-GS1-ALM	SDM-GS5-ALM
1	1.00	1.00	1.00	1.00
8	5.54	5.23	4.38	4.78
16	8.89	8.55	6.61	7.07
32	11.69	11.94	8.19	8.89
Lagr. value	− 127.37	− 127.37	− 127.71	− 127.58

No. Proc.	Speedup for SSLP 10-50-500			
	OOQP	PIPS-IPM	SDM-GS1-ALM	SDM-GS5-ALM
1	1.00	1.00	1.00	1.00
8	2.64	2.80	6.87	6.95
16	2.70	2.92	12.95	12.84
32	2.98	3.40	21.67	20.98
Lagr. value	− 349.14	− 349.14	− 349.48	− 349.14

No. Proc.	Speedup for SSLP 10-50-2000	
	SDM-GS1-ALM	SDM-GS5-ALM
1	1.00	1.00
2	2.34	2.34
4	4.81	4.83
8	9.29	9.25
16	18.69	18.48
32	34.63	35.10
64	60.59	60.93
Lagr. value	− 348.35	− 347.75

We draw the following conclusions from the results of the Test 2 experiments reported in Tables 1 and 2.

1. The improvement in parallel speedup (SDM-GS-ALM columns) over either OOQP or PIPS-IPM is evident for all problems except for the one with the fewest number of scenarios (SSLP 5-25-100).
2. Slightly inferior final Lagrange bounds reported for SDM-GS1-ALM ($t_{max} = 1$) are evident. This deficit is improved by using SDM-GS with $t_{max} = 5$, as done for the SDM-GS5-ALM experiments. But even these bounds are usually not as good as the bounds obtained with OOQP or PIPS-IPM; this is due to their more exact solving of the master problem instances. This suggests that as the iterations $k \geq 1$ increase, it is advantageous to solve the continuous master problem with SDM-GS iterations using larger t_{max} values.

Table 2 DCAP: comparing speedup and final best Lagrangian bound

No. Proc.	Speedup for DCAP 233-500			
	OOQP	PIPS-IPM	SDM-GS1-ALM	SDM-GS5-ALM
1	1.00	1.00	1.00	1.00
8	2.44	5.32	6.88	8.11
16	2.81	8.15	13.28	15.65
32	1.63	10.25	23.42	27.40
Lagr. value	1736.68	1736.68	1734.99	1736.02

No. Proc.	Speedup for DCAP 243-500			
	OOQP	PIPS-IPM	SDM-GS1-ALM	SDM-GS5-ALM
1	1.00	1.00	1.00	1.00
8	2.85	5.71	6.51	7.61
16	3.59	5.85	12.28	14.44
32	1.98	6.44	21.99	25.25
Lagr. value	2165.48	2165.50	2162.58	2164.48

No. Proc.	Speedup for DCAP 332-500			
	OOQP	PIPS-IPM	SDM-GS1-ALM	SDM-GS5-ALM
1	1.00	1.00	1.00	1.00
8	2.03	5.56	6.83	8.50
16	2.33	5.00	12.84	16.20
32	1.21	6.61	21.83	23.48
Lagr. value	1587.44	1587.44	1584.77	1586.11

No. Proc.	Speedup for DCAP 342-500			
	OOQP	PIPS-IPM	SDM-GS1-ALM	SDM-GS5-ALM
1	1.00	1.00	1.00	1.00
8	2.45	3.78	7.16	8.25
16	2.71	4.36	12.95	15.49
32	1.84	4.64	22.41	26.93
Lagr. value	1902.84	1903.21	1900.81	1901.90

3. Interestingly, parallel speedup is enhanced for SDM-GS5-ALM over SDM-GS1-ALM; although the latter yields lower average total wall clock time per iteration, the proportion of efficiently parallelizable work seems to increase in the former.

For Test 2, we also tested the performance of Algorithm 3 on the SSLP 10-50-2000 problem, which is of substantially larger scale than the other test problems considered in this paper. Using $N = 1, 2, 4, 8, 16, 32, 64$ processors, we see very good speedup, which suggests the realized benefit of distributing the use of memory. We also see that for such large-scale problems, the additional cost in time of performing more inner

loop Gauss–Seidel iterations (larger t_{max}) becomes marginal, since the cost of solving the mixed-integer linear subproblems takes a larger share of the computational time.

6 Conclusion and future work

Our contribution is motivated by the goal of improving the efficiency of parallelization applied to iterative approaches for solving the Lagrangian dual problem of large scale optimization problems. These problems have nonlinear convex differentiable objective f , decomposable nonconvex constraint set X , and nondecomposable affine constraint set $Qx = z$ to which Lagrangian relaxation is applied. Problems of such a form include the split variable extensive form of mixed-integer linear stochastic programs as a special case. Implicitly, our approach refers to the convex hull $\text{conv}(X)$ of X , and the assumed lack of known description of $\text{conv}(X)$ needs to be addressed. Proximal bundle methods (alternatively in the form of the proximal simplicial decomposition method or stabilized column generation) are well-known for addressing the latter issue. In the former issue, that of exploiting the large scale structure to apply parallel computation efficiently, we develop a modified augmented Lagrangian (AL) method with approximate subproblem solutions that incorporates ideas from the proximal bundle method.

The approximation of subproblem solutions is based on an iterative approach that integrates ideas from the simplicial decomposition method (SDM) (for constructing inner approximations of $\text{conv}(X)$) and the nonlinear block Gauss–Seidel method. It is the latter Gauss–Seidel aspect that is primarily responsible for enhancing the parallel efficiency that is observed in the numerical experiments. While convergence analysis of the integrated SDM-GS approach may be derived from slight modifications to results in [21], for the sake of completeness and explicitness, we provide in the appendix a proof of optimal convergence of SDM-GS as it is applied within our algorithm under a standard set of conditions. A distinction between so-called “serious” steps and “null” steps, in analogy to the proximal bundle method, is also recovered. Once these aspects are successfully integrated, then the contribution is complete, where the beneficial stabilization associated with proximal point methods and the ability to apply parallelization more efficiently are both realized. The resulting algorithm developed in this paper is referred to as SDM-GS-ALM, which has similar functionality to the alternating direction method of multipliers (ADMM).

We performed numerical tests of two sorts. In Test 1, we examined the impact of varying the serious step condition parameter. We found that parameterizations that effect more stringent serious step conditions seem to have the effect of mitigating the early iteration instability due to penalty parameters that are too large. At the same time, the more stringent serious step condition parameterizations seemed to result in slower convergence to dual optimality in the tail-end. As is the case for proximal bundle methods, information obtained in the serious step condition tests may be used to beneficially adjust the proximal term penalty coefficient in early iterations.

In Test 2, we examined the efficiency of parallelization, measured by the speedup ratio, due to the use of the SDM-GS-ALM, compared versus pre-existing implementations of the proximal bundle method that use structure exploiting primal dual interior

point methods to improve parallel efficiency. We saw in these results a promising increase in parallel efficiency due to the use of SDM-GS-ALM, where the increase in parallel efficiency is attributed primarily to the successful incorporation of Gauss–Seidel iterations. The results of the last problem tested, SSLP 10-50-2000, additionally suggested a benefit due to the ability of SDM-GS-ALM to distribute not just the workload, but also the use of memory. The vector of auxiliary variables z is the only substantial block of data that needs to be stored and modified by all processors. In the context of stochastic optimization problems, this represents a modest communication bottleneck in proportion to the number of first-stage variables for two-stage problems, while for multistage problems, the amount of such data that must be stored by every processor and modified by parallel communication can increase exponentially with the number of stages.

Potential future improvements include the following. While a default implementation of SDM-GS-ALM would have one Gauss–Seidel iteration per SDM-GS call, the Lagrangian bounds reported from the Test 2 experiments suggest that an improved implementation would have early iterations use one Gauss–Seidel iteration per SDM-GS call, but steadily increase the number of Gauss–Seidel iterations per SDM-GS call for the later iterations. This results in better Lagrangian bounds at termination. While these extra Gauss–Seidel iterations require extra parallel communication, the additional wall clock time required becomes increasingly marginal for larger problems where the cost of solving the SDM linearized subproblems associated with expanding the inner approximation increasingly outweighs the cost associated with computing the approximate solution of the continuous master problem and any required parallel communications.

A potentially large improvement to the speed of convergence, in terms of wall clock time, would be to incorporate into the analysis the degree to which the SDM linearized subproblem can be solved suboptimally and yet retain the optimal convergence. We expect that solving these subproblems exactly, particularly in the early iterations, is highly wasteful, and providing a theoretical basis for controlling the tolerance of solution inaccuracy would be of great value. Another potential avenue for future work is to extend the experimental analysis to multistage mixed-integer stochastic optimization problems and/or nonlinear problems, as the form of the problem addressed by SDM-GS-ALM is general enough to model these types of problems.

A Technical lemmas for establishing optimal convergence of SDM-GS

Given initial $(x^0, z^0) \in X \times Z \subset \mathbb{R}^n \times \mathbb{R}^q$, we consider the generation of the sequence $\{(x^k, z^k)\}$ with iterations computed using Algorithm 4, whose target problem is given by

$$\min_{x, z} \{F(x, z) : x \in X, z \in Z\}, \quad (32)$$

where $(x, z) \mapsto F(x, z)$ is convex and continuously differentiable over $X \times Z$, and sets X and Z are closed and convex, with X bounded and $z \mapsto F(x, z)$ is inf-compact for each $x \in X$.

We define the directional derivative with respect to x as

$$F'_x(x, z; d) := \lim_{\alpha \downarrow 0} \frac{F(x + \alpha d, z) - F(x, z)}{\alpha}.$$

Of interest is the satisfaction of the following local stationarity condition at $x \in X$:

$$F'_x(x, z; d) \geq 0 \quad \text{for all } d \in X - \{x\} \quad (33)$$

for any limit point $(x, z) = (\bar{x}, \bar{z})$ of some sequence $\{(x^k, z^k)\}$ of feasible solutions to problem (32). For the sake of nontriviality, we shall assume that the x -stationarity condition (33) never holds at $(x, z) = (x^k, z^k)$ for any $k \geq 0$. Thus, for each $x^k, k \geq 0$, there always exists a $d^k \in X - \{x^k\}$ for which $F'_x(x^k, z^k; d^k) < 0$.

Direction Assumptions (DAs) For each iteration $k \geq 0$, given $x^k \in X$ and $z^k \in Z$, we have d^k chosen so that (1) $x^k + d^k \in X$; and (2) $F'_x(x^k, z^k; d^k) < 0$.

Gradient Related Assumption (GRA) Given a sequence $\{(x^k, z^k)\}$ with $\lim_{k \rightarrow \infty} (x^k, z^k) = (\bar{x}, \bar{z})$, and a bounded sequence $\{d^k\}$ of directions, then the existence of a direction $\bar{d} \in X - \{\bar{x}\}$ such that $F'_x(\bar{x}, \bar{z}; \bar{d}) < 0$ implies that

$$\limsup_{k \rightarrow \infty} F'_x(x^k, z^k; d^k) < 0. \quad (34)$$

In this case, we say that $\{d^k\}$ is *gradient related* to $\{x^k\}$. This gradient related condition is similar to the one defined in [3]. The sequence of directions d^k is typically gradient related to $\{x^k\}$ by construction. (See Lemma 7.)

To state the last assumption, we require the notion of an Armijo rule step length $\alpha^k \in (0, 1]$ given (x^k, z^k, d^k) and parameters $\beta, \sigma \in (0, 1)$.

Algorithm 4 Computing an Armijo rule step length α^k at iteration k .

```

1: function ARMIJOSTEP( $F, x^k, z^k, d^k, \beta, \sigma$ )
2:    $\alpha^k \leftarrow 1$ 
3:   while  $F(x^k + \alpha^k d^k, z^k) - F(x^k, z^k) > \alpha^k \sigma F'_x(x^k, z^k; d^k)$  do
4:      $\alpha^k \leftarrow \beta \alpha^k$ 
5:   end while
6:   return  $\alpha^k$ 
7: end function

```

Remark 6 Under mild assumptions on F such as continuity that guarantee the existence of finite $F'_x(x, z; d)$ for all $(x, z, d) \in \{(x, z, d) : x \in X, d \in X - \{x\}, z \in Z\}$, we may assume that the while loop of Lines 3–5 terminates after a finite number of iterations. Thus, we have $\alpha^k \in (0, 1]$ for each $k \geq 1$.

The last significant assumption is stated as follows.

Sufficient Decrease Assumption (SDA) For sequences $\{(x^k, z^k, d^k)\}$ and step lengths $\{\alpha^k\}$ computed according to Algorithm 4, we assume for each $k \geq 0$, that (x^{k+1}, z^{k+1}) satisfies

$$F(x^{k+1}, z^{k+1}) \leq F(x^k + \alpha^k d^k, z^k).$$

Lemma 6 For problem (32), let $F : \mathbb{R}^{n_x} \times \mathbb{R}^{n_z} \mapsto \mathbb{R}$ be convex and continuously differentiable, $X \subset \mathbb{R}^{n_x}$ convex and compact, and $Z \subseteq \mathbb{R}^{n_z}$ closed and convex. Furthermore, assume for each $x \in X$ that $z \mapsto F(x, z)$ is inf-compact. If a sequence $\{(x^k, z^k, d^k)\}$ satisfies the DA, the GRA, and the SDA for some fixed $\beta, \sigma \in (0, 1)$, then the sequence (x^k, z^k) has limit points (\bar{x}, \bar{z}) , each of which satisfies the stationarity condition (33).

Proof The existence of limit points (\bar{x}, \bar{z}) follows from the compactness of X , the inf-compactness of $z \mapsto F(x, z)$ for each $x \in X$, and the SDA. In generating $\{\alpha^k\}$ according to the Armijo rule as implemented in Lines 2–5 of Algorithm 4, we have

$$\frac{F(x^k + \alpha^k d^k, z^k) - F(x^k, z^k)}{\alpha^k} \leq \sigma F'_x(x^k, z^k; d^k). \quad (35)$$

By the DA, $F'_x(x^k, z^k; d^k) < 0$ and since $\alpha^k > 0$ for each $k \geq 1$ by Remark 6, we infer from (35) that $F(x^k + \alpha^k d^k, z^k) < F(x^k, z^k)$. By construction, we have $F(x^{k+1}, z^{k+1}) \leq F(x^k + \alpha^k d^k, z^k) < F(x^k, z^k)$. By the monotonicity $F(x^{k+1}, z^{k+1}) < F(x^k, z^k)$ and F being bounded from below on $X \times Z$, we have $\lim_{k \rightarrow \infty} F(x^k, z^k) = \bar{F} > -\infty$. Therefore,

$$\lim_{k \rightarrow \infty} F(x^{k+1}, z^{k+1}) - F(x^k, z^k) = 0,$$

which implies

$$\lim_{k \rightarrow \infty} F(x^k + \alpha^k d^k, z^k) - F(x^k, z^k) = 0. \quad (36)$$

We assume for sake of contradiction that $\lim_{k \rightarrow \infty} (x^k, z^k) = (\bar{x}, \bar{y})$ does not satisfy the stationarity condition (33). By GRA, we have that $\{d^k\}$ is gradient related to $\{x^k\}$; that is,

$$\limsup_{k \rightarrow \infty} F'_x(x^k, z^k; d^k) < 0. \quad (37)$$

Thus, it follows from (35)–(37) that $\lim_{k \rightarrow \infty} \alpha^k = 0$.

Consequently, after a certain iteration $k \geq \bar{k}$, we can define $\{\bar{\alpha}^k\}$, $\bar{\alpha}^k = \alpha^k / \beta$, where $\bar{\alpha}^k \leq 1$ for $k \geq \bar{k}$, and so we have

$$\sigma F'_x(x^k, z^k; d^k) < \frac{F(x^k + \bar{\alpha}^k d^k, z^k) - F(x^k, z^k)}{\bar{\alpha}^k}. \quad (38)$$

Since F is continuously differentiable, the mean value theorem may be applied to the right-hand side of (38) to get

$$\sigma F'_x(x^k, z^k; d^k) < F'_x(x^k + \tilde{\alpha}^k d^k, z^k; d^k), \quad (39)$$

for some $\tilde{\alpha}^k \in [0, \bar{\alpha}^k]$.

Again, using the assumption $\limsup_{k \rightarrow \infty} F'_x(x^k, z^k; d^k) < 0$, and also the compactness of $X - X$, we take a limit point \bar{d} of $\{d^k\}$, with its associated subsequence index set denoted by \mathcal{K} , such that $F'_x(\bar{x}, \bar{z}, \bar{d}) < 0$. Taking the limits over the subsequence indexed by \mathcal{K} , we have $\lim_{k \rightarrow \infty, k \in \mathcal{K}} F'_x(x^k, z^k; d^k) = F'_x(\bar{x}, \bar{z}, \bar{d})$ and $\lim_{k \rightarrow \infty, k \in \mathcal{K}} F'_x(x^k + \tilde{\alpha}^k d^k, z^k; d^k) = F'_x(\bar{x}, \bar{z}, \bar{d})$. These two limits holds since (1) $(x, z) \mapsto F'_x(x, z; d)$ for each $d \in X - X$ is continuous and (2) $d \mapsto F'_x(x, z; d)$ is locally Lipschitz continuous for each $(x, z) \in X \times Z$ (e.g., Proposition 2.1.1 of [60]); these two facts together imply that $(x, z; d) \mapsto F'_x(x, z; d)$ is continuous. Then, inequality (39) becomes in the limit as $k \rightarrow \infty, k \in \mathcal{K}$,

$$\sigma F'_x(\bar{x}, \bar{z}, \bar{d}) \leq F'_x(\bar{x}, \bar{z}, \bar{d}) \implies 0 \leq (1 - \sigma) F'_x(\bar{x}, \bar{z}, \bar{d}).$$

Since $(1 - \sigma) > 0$ and $F'_x(\bar{x}, \bar{z}, \bar{d}) < 0$, we have a contradiction. Thus, \bar{x} must satisfy the stationary condition (33). \square

Remark 7 Noting that $F'_x(x^k, z^k; d^k) = \nabla_x F(x^k, z^k)^\top d^k$ under the assumption of continuous differentiability of F , one means of constructing $\{d^k\}$ is as follows:

$$d^k \leftarrow \operatorname{argmin}_d \left\{ \nabla_x F(x^k, z^k)^\top d : d \in X - \{x^k\} \right\}. \quad (40)$$

Lemma 7 Given sequence $\{(x^k, z^k)\}$ with $\lim_{k \rightarrow \infty} (x^k, z^k) = (\bar{x}, \bar{z})$, let each d^k , $k \geq 1$, be generated as in (40). Then $\{d^k\}$ is gradient related to $\{x^k\}$.

Proof By the construction of d^k , $k \geq 1$, we have

$$F'_x(x^k, z^k; d^k) \leq F'_x(x^k, z^k; d) \quad \forall d \in X - \{x^k\}.$$

Taking the limit, we have

$$\limsup_{k \rightarrow \infty} F'_x(x^k, z^k; d^k) \leq \limsup_{k \rightarrow \infty} F'_x(x^k, z^k; d) \leq F'_x(\bar{x}, \bar{z}; d) \quad \forall d \in X - \{\bar{x}\},$$

where the last inequality follows from the upper semicontinuity of the function $(x, z, d) \mapsto F'_x(x, z; d)$, which holds in our setting due, primarily, to Proposition 2.1.1 (b) of [60] given that F is assumed to be convex and continuous on \mathbb{R}^n . Taking

$$\bar{d} \in \operatorname{argmin}_d \left\{ F'_x(\bar{x}, \bar{z}; d) : d \in X - \{\bar{x}\} \right\},$$

we have by the assumed nonstationarity that $F'_x(\bar{x}, \bar{z}, \bar{d}) < 0$. Thus, $\limsup_{k \rightarrow \infty} F'_x(x^k, z^k; d^k) < 0$, and so GRA holds. \square

References

1. Birge, J.R., Louveaux, F.: *Introduction to Stochastic Programming*. Springer, Berlin (2011)
2. Shor, N.: *Minimization Methods for Non-differentiable Functions*. Springer, New York (1985)
3. Bertsekas, D.: *Nonlinear Programming*. Athena Scientific, Belmont (1999)
4. Ruszczyński, A.: *Nonlinear Optimization*. Princeton University Press, Princeton (2006)
5. Carøe, C.C., Schultz, R.: Dual decomposition in stochastic integer programming. *Oper. Res. Lett.* **24**(1), 37–45 (1999)
6. Bertsekas, D.: *Constrained Optimization and Lagrange Multiplier Methods*. Academic Press, London (1982)
7. Hestenes, M.R.: Multiplier and gradient methods. *J. Optim. Theory Appl.* **4**, 303–320 (1969)
8. Powell, M.J.D.: A method for nonlinear constraints in minimization problems. In: Fletcher, R. (ed.) *Optimization*. Academic Press, New York (1969)
9. Rockafellar, R.: Monotone operators and the proximal point algorithm. *SIAM J. Control Optim.* **14**(5), 877–898 (1976)
10. Lemaréchal, C.: An Extension of Davidson Methods to Non differentiable Problems, pp. 95–109. Springer, Berlin Heidelberg (1975)
11. de Oliveira, W., Sagastizábal, C.: Bundle methods in the XX1st century: a bird's-eye view. *Pesqui. Oper.* **34**, 647–670 (2014)
12. Hare, W., Sagastizábal, C., Solodov, M.: A proximal bundle method for nonsmooth nonconvex functions with inexact information. *Comput. Optim. Appl.* **63**, 1–28 (2016)
13. de Oliveira, W., Sagastizábal, C., Lemaréchal, C.: Convex proximal bundle methods in depth: a unified analysis for inexact oracles. *Math. Program.* **148**(1), 241–277 (2014)
14. Amor, H.B., Desrosiers, J., Frangioni, A.: On the choice of explicit stabilizing terms in column generation. *Discrete Appl. Math.* **157**(6), 1167–1184 (2009)
15. Bertsekas, D.: Incremental aggregated proximal and augmented Lagrangian algorithms. *arXiv preprint arXiv:1509.09257* (2015)
16. Fischer, F., Helmberg, C.: A parallel bundle framework for asynchronous subspace optimization of nonsmooth convex functions. *SIAM J. Optim.* **24**(2), 795–822 (2014)
17. Lubin, M., Martin, K., Petra, C., Sandıkçı, B.: On parallelizing dual decomposition in stochastic integer programming. *Oper. Res. Lett.* **41**(3), 252–258 (2013)
18. Bertsekas, D.: *Convex Optimization Algorithms*. Athena Scientific, Belmont (2015)
19. Holloway, C.: An extension of the Frank and Wolfe method of feasible directions. *Math. Program.* **6**(1), 14–27 (1974)
20. Von Hohenbalken, B.: Simplicial decomposition in nonlinear programming algorithms. *Math. Program.* **13**(1), 49–68 (1977)
21. Bonettini, S.: Inexact block coordinate descent methods with application to non-negative matrix factorization. *IMA J. Numer. Anal.* **31**(4), 1431–1452 (2011)
22. Grippo, L., Sciandrone, M.: On the convergence of the block nonlinear Gauss–Seidel method under convex constraints. *Oper. Res. Lett.* **26**(3), 127–136 (2000)
23. Hildreth, C.: A quadratic programming procedure. *Nav. Res. Logist. Q.* **4**, 79–85 (1957). 361
24. Tseng, P.: Convergence of a block coordinate descent method for nondifferentiable minimization. *J. Optim. Theory Appl.* **109**, 475–494 (2001)
25. Warga, J.: Minimizing certain convex functions. *SIAM J. Appl. Math.* **11**, 588–593 (1963)
26. Eckstein, J.: A practical general approximation criterion for methods of multipliers based on Bregman distances. *Math. Program.* **96**(1), 61–86 (2003)
27. Eckstein, J., Silva, P.: A practical relative error criterion for augmented lagrangians. *Math. Program.* **141**(1), 319–348 (2013)
28. Hamdi, A., Mahey, P., Dussault, J.P.: Recent Advances in Optimization: Proceedings of the 8th French–German Conference on Optimization Trier, July 21–26, 1996, chap. A New Decomposition Method in Nonconvex Programming via a Separable Augmented Lagrangian, pp. 90–104. Springer Berlin Heidelberg, Berlin, Heidelberg (1997)
29. Boyd, S., Parikh, N., Chu, E., Peleato, B., Eckstein, J.: Distributed optimization and statistical learning via the alternating direction method of multipliers. *Found. Trends Mach. Learn.* **3**(1), 1–122 (2011)
30. Gabay, D., Mercier, B.: A dual algorithm for the solution of nonlinear variational problems via finite element approximation. *Comput. Math. Appl.* **2**, 17–40 (1976)

31. Glowinski, R., Marrocco, A.: Sur l'approximation, par elements finis d'ordre un, et la resolution, par penalisation-dualité, d'une classe de problems de dirichlet non lineares. *Revue Française d'Automatique, Informatique, et Recherche Opérationnelle* **9**, 41–76 (1975)
32. Chatzipanagiotis, N., Dentcheva, D., Zavlanos, M.: An augmented Lagrangian method for distributed optimization. *Math. Program.* **152**(1), 405–434 (2014)
33. Tappenden, R., Richtárik, P., Büke, B.: Separable approximations and decomposition methods for the augmented Lagrangian. *Optim. Methods Softw.* **30**(3), 643–668 (2015)
34. Mulvey, J., Ruszczyński, A.: A diagonal quadratic approximation method for large scale linear programs. *Oper. Res. Lett.* **12**(4), 205–215 (1992)
35. Ruszczyński, A.: On convergence of an augmented Lagrangian decomposition method for sparse convex optimization. *Math. Oper. Res.* **20**(3), 634–656 (1995)
36. Kiwiel, K., Rosa, C., Ruszczyński, A.: Proximal decomposition via alternating linearization. *SIAM J. Optim.* **9**(3), 668–689 (1999)
37. Lin, X., Pham, M., Ruszczyński, A.: Alternating linearization for structured regularization problems. *J. Mach. Learn. Res.* **15**, 3447–3481 (2014)
38. Chen, G., Teboulle, M.: A proximal-based decomposition method for convex minimization problems. *Math. Program.* **64**, 81–101 (1994)
39. He, B., Liao, L.Z., Han, D., Yang, H.: A new inexact alternating directions method for monotone variational inequalities. *Math. Program.* **92**, 103–118 (2002)
40. Boland, N., Christiansen, J., Dandurand, B., Eberhard, A., Linderoth, J., Luedtke, J., Oliveira, F.: Progressive hedging with a Frank–Wolfe based method for computing stochastic mixed-integer programming Lagrangian dual bounds. *Optimization Online* (2016). http://www.optimization-online.org/DB_HTML/2016/03/5391.html
41. Eckstein, J., Bertsekas, D.: On the Douglas–Rachford splitting method and the proximal point algorithm for maximal monotone operators. *Math. Program.* **55**(1–3), 293–318 (1992)
42. Eckstein, J., Yao, W.: Understanding the Convergence of the Alternating Direction Method of Multipliers: Theoretical and Computational Perspectives. Rutgers University, New Brunswick (2014). Tech. rep
43. Mahey, P., Oualibouch, S., Tao, P.D.: Proximal decomposition on the graph of a maximal monotone operator. *SIAM J. Optim.* **5**(2), 454–466 (1995)
44. Feizollahi, M.J., Costley, M., Ahmed, S., Grijalva, S.: Large-scale decentralized unit commitment. *Int. J. Electr. Power Energy Syst.* **73**, 97–106 (2015)
45. Gade, D., Hackebeil, G., Ryan, S.M., Watson, J.P., Wets, R.J.B., Woodruff, D.L.: Obtaining lower bounds from the progressive hedging algorithm for stochastic mixed-integer programs. *Math. Program.* **157**(1), 47–67 (2016)
46. Rockafellar, R.T.: Augmented Lagrangians and applications of the proximal point algorithm in convex programming. *Math. Oper. Res.* **1**(2), 97–116 (1976)
47. Rockafellar, R.: *Convex Analysis*. Princeton University Press, Princeton (1970)
48. Hathaway, R.J., Bezdek, J.C.: Grouped coordinate minimization using Newton's method for inexact minimization in one vector coordinate. *J. Optim. Theory Appl.* **71**(3), 503–516 (1991)
49. Kiwiel, K.C.: Approximations in proximal bundle methods and decomposition of convex programs. *J. Optim. Theory Appl.* **84**(3), 529–548 (1995)
50. Bodur, M., Dash, S., Günlük, O., Luedtke, J.: Strengthened Benders cuts for stochastic integer programs with continuous recourse (2014). http://www.optimization-online.org/DB_FILE/2014/03/4263.pdf. Last Accessed 13 Jan (2015)
51. Ahmed, S., Garcia, R., Kong, N., Ntamo, L., Parija, G., Qiu, F., Sen, S.: SIPLIB: A stochastic integer programming test problem library (2015). <http://www.isye.gatech.edu/sahmed/siplib>
52. Ntamo, L.: Decomposition algorithms for stochastic combinatorial optimization: Computational experiments and extensions. Ph.D. thesis (2004)
53. The MathWorks, Natick: MATLAB 2012b (2014)
54. IBM Corporation: IBM ILOG CPLEX Optimization Studio CPLEX Users Manual. http://www.ibm.com/support/knowledgecenter/en/SSSA5P_12.6.1/ilog.odms.studio.help/pdf/usrcplex.pdf. Last Accessed 22 Aug (2016)
55. IBM Corporation: IBM ILOG CPLEX V12.5. <http://www-01.ibm.com/software/commerce/optimization/cplex-optimizer/>. Last Accessed 28 Jan (2016)
56. COmputational INfrastructure for Operations Research. <http://www.coin-or.org/>. Last Accessed 28 Jan (2016)

57. National Computing Infrastructure (NCI): NCI Website. <http://www.nci.org.au>. Last Accessed 19 Nov 2016
58. Gertz, E., Wright, S.: Object-oriented software for quadratic programming. *ACM Trans. Math. Softw.* **29**(1), 58–81 (2003)
59. Lubin, M., Petra, C., Anitescu, M., Zavala, V.: Scalable stochastic optimization of complex energy systems. In: *Proceedings of 2011 International Conference for High Performance Computing, Networking, Storage and Analysis*, pp. 64:164:10. ACM, Seattle, WA (2011)
60. Clarke, F.: *Optimization and Nonsmooth Analysis*. Society for Industrial and Applied Mathematics (1990)