



An adaptive primal-dual framework for nonsmooth convex minimization

Quoc Tran-Dinh¹ · Ahmet Alacaoglu² · Olivier Fercoq³ · Volkan Cevher²

Received: 21 June 2018 / Accepted: 19 September 2019 / Published online: 31 October 2019
© Springer-Verlag GmbH Germany, part of Springer Nature and Mathematical Optimization Society 2019

Abstract

We propose a new self-adaptive and double-loop smoothing algorithm to solve composite, nonsmooth, and constrained convex optimization problems. Our algorithm is based on Nesterov's smoothing technique via general Bregman distance functions. It self-adaptively selects the number of iterations in the inner loop to achieve a desired complexity bound without requiring to set the accuracy a priori as in variants of augmented Lagrangian methods (ALM). We prove $\mathcal{O}(\frac{1}{k})$ -convergence rate on the last iterate of the outer sequence for both unconstrained and constrained settings in contrast to ergodic rates which are common in ALM as well as alternating direction method-of-multipliers literature. Compared to existing inexact ALM or quadratic penalty methods, our analysis does not rely on the worst-case bounds of the subproblem solved by the inner loop. Therefore, our algorithm can be viewed as a restarting technique applied to the ASGARD method in Tran-Dinh et al. (SIAM J Optim 28(1):96–134, 2018) but with rigorous theoretical guarantees or as an inexact ALM with explicit inner loop termination rules and adaptive parameters. Our algorithm only requires to initialize the parameters once, and automatically updates them during the iteration process without tuning. We illustrate the superiority of our methods via several examples as compared to the state-of-the-art.

Keywords Primal-dual first-order methods · Restarting · Augmented Lagrangian · Self-adaptive method · Nonsmooth convex optimization

Mathematics Subject Classification 90C25 · 90C06 · 90-08

1 Introduction

Problem settings: We study the following nonsmooth composite convex minimization template:

$$P^* := \min_{x \in \mathbb{R}^p} \left\{ P(x) := f(x) + g(Ax) \right\}, \quad (1)$$

Extended author information available on the last page of the article

where both $f : \mathbb{R}^p \rightarrow \mathbb{R} \cup \{+\infty\}$ and $g : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ are proper, closed, and nonsmooth convex functions, and $A : \mathbb{R}^p \rightarrow \mathbb{R}^n$ is a linear operator. In addition, we assume that g is Lipschitz continuous when restricted to its domain (see Assumption 1). This assumption covers many important cases:

1. g is Lipschitz continuous on \mathbb{R}^n , i.e. there exists $M_g \in (0, +\infty)$ such that $|g(x) - g(y)| \leq M_g \|x - y\|$ for any $x, y \in \mathbb{R}^n$.
2. $g := \delta_{b+\mathcal{K}}$, the indicator of $b + \mathcal{K}$, where $b \in \mathbb{R}^n$ is given and \mathcal{K} is a given nonempty, closed, and convex set in \mathbb{R}^n . In this case, (1) automatically covers the following constrained convex problem:

$$f^* := \min_{x \in \mathbb{R}^p} \left\{ f(x) \text{ s.t. } Ax - b \in \mathcal{K} \right\}. \tag{2}$$

In particular, if $\mathcal{K} = \{0^n\}$, then $Ax - b \in \mathcal{K}$ reduces to $Ax = b$.

3. g is a polyhedral function, i.e. its epigraph is a polyhedron.

Our goals: Our goal is to design new adaptive primal-dual methods to solve both (1) and (2) that have low per-iteration complexity cost, i.e., they only require the proximal operators of f and g , and matrix-vector multiplications Ax and $A^T y$, while having the best-known non-averaging convergence rates. We require that both problems are convex and satisfy strong duality assumptions. We do not assume that the problems are strongly convex or smooth.

Composite versus constrained settings: For the general setting (1), under different choices of f and g , it covers a wide range of applications from different fields including compressive sensing, image and signal processing, machine learning, statistics, operations research, and optimal control. Classical and well-known examples such as LASSO, square-root LASSO, support vector machine, image denoising and deblurring, and matrix completion can be cast into (1), see, e.g., [10,19,52,65] for some concrete examples.

For the setting (2), we do not impose any restriction on \mathcal{K} . Hence, it covers a large class of constrained problems including equality and inequality constraints. When \mathcal{K} is a given cone (e.g., \mathbb{R}_+^n , second-order cone, or symmetric positive semidefinite cone), problem (2) also covers problems with cone constraints such as linear programming, second-order cone, and semidefinite programming. Although the theory for (1) as well as for (2) are well developed, various numerical methods for solving these problems rely on different structure assumptions and do not have a unified analysis: cf., Sect. 5.

Related works: Under only convexity and zero duality gap assumptions, the state-of-the-art methods for solving (1) include first-order primal-dual methods (FOPDM) [2,13], and augmented Lagrangian-based algorithms [5,35,55]. While FOPDM directly tackles problem (1), the augmented Lagrangian-based framework (ALM) and its variants solve (1) via a constrained reformulation as follows (or using other forms):

$$P^* := \min_{x \in \mathbb{R}^p, z \in \mathbb{R}^n} \left\{ P(x, z) := f(x) + g(z) \text{ s.t. } Ax - z = 0 \right\}. \tag{3}$$

Alternating direction method of multipliers (ADMM) proposed in [29,33] is another (and perhaps the most) successful method to solve (3). One can derive a variant of

ADMM by applying an alternating strategy between x and z to the ALM framework to break the computational bottleneck in the primal subproblem, see, e.g. [10]. Inexact and linearized variants enhance the scalability of ALM and ADMM for the same problem template [51,66,67].

While ADMM and FOPDM and their variants work really well in practice, their best-known convergence rate is $\mathcal{O}\left(\frac{1}{k}\right)$ under only convexity and zero duality gap assumptions, where k is the iteration counter. Moreover, such a rate is achieved via an ergodic sense (i.e., using an averaging sequence or a weighted averaging sequence) [13,14,21,22,41,42,56].

In stark contrast, empirical evidence shows that averaging sequences in FOPDM and ALM exhibit the theoretical worst case rate $\mathcal{O}\left(\frac{1}{k}\right)$ in practice compared to the last iterate of the algorithm (see Sect. 4.1 for a concrete example), which is superior and often locally linear in many examples.¹ However, for these methods, the last iterate sequence generally has asymptotic or non-asymptotic convergence guarantee, but it is much slower than that of averaging sequences, e.g., $\mathcal{O}\left(\frac{1}{\sqrt{k}}\right)$ -rate [22].

Recently, [59] proposed an *accelerated smoothed gap reduction* (ASGARD) framework to solve nonsmooth convex optimization problems. ASGARD combines acceleration, smoothing, and homotopy techniques to handle both unconstrained and constrained nonsmooth problems, including (1) and (2).

One notable feature of ASGARD is a non-ergodic optimal $\mathcal{O}\left(\frac{1}{k}\right)$ rate on the objective residual, and feasibility violation in the constrained settings. Moreover, this method only requires one proximal operator of f , one matrix-vector multiplication, and one adjoint operator per iteration. When f is separable, the algorithm can be naturally parallelized. However, as also noted in [59], ASGARD needs restarting to be competitive with state-of-the-art methods such as ADMM and FOPDM in practice. This is not surprising since empirical evidence [28,32,50,57] has shown that restarting significantly improves the actual convergence rate in practice. While there exists theory to support the restarting strategies in accelerated gradient-type methods, supporting theory of these strategies are not yet investigated in primal-dual methods.

Contributions: In this paper, we introduce an analysis framework for restarting ASGARD and prove the same worst-case $\mathcal{O}\left(\frac{1}{k}\right)$ rate in a non-ergodic sense. While doing so, we identify that restarting ASGARD corresponds to an inexact ALM algorithm in the constrained case. In contrast to existing works on this front, our method has explicit inner-loop termination rules and does not need to set a horizon (i.e., the maximum number of inner iterations or a predefined inner loop accuracy) for the algorithm. As a result, we present a method which has the guarantees on the last iterate compared to ALM/ADMM methods and extend the guarantees of ASGARD to the restarting case which significantly improves the practical performance. In addition, we allow general Bregman distances to be used for smoothing and proximal operators in contrast to the original ASGARD scheme. A more thorough discussion and comparison between our method and state-of-the-arts is deferred to Sect. 5. To this end, our contribution can be summarized as follows.

- (a) We propose a new self-adaptive, double-loop smoothing algorithm to solve non-smooth convex optimization problems of the form (1). Our algorithm is based

¹ There exist examples showing arbitrarily slow convergence rate of ADMM, see, e.g., [22].

on Nesterov's smoothing technique via general Bregman distance functions. It self-adaptively selects the number of iterations in the inner loop to achieve a desired complexity bound without requiring the accuracy a priori as in variants of ALM. Compared to ASGARD [59], it incorporates restarts, updates the dual center, and can work with general Bregman distances instead of only Lipschitz gradient distances.

- (b) We prove $\mathcal{O}\left(\frac{1}{k}\right)$ -convergence rate on the last iterate of the outer sequence for both unconstrained and constrained settings in contrast to ergodic rates which are commonly seen in ALM/ADMM literature. This rate is known to be optimal [46,48,64] under just convexity and strong duality assumptions and when k is in the order of p . Compared to existing inexact ALM or quadratic penalty methods such as [43,68], our analysis does not rely on the worst-case bounds of the subproblem solved by the inner loop. Therefore, our algorithm can be viewed as a restarting technique applied to ASGARD but with rigorous theoretical guarantees or as an inexact ALM with explicit inner loop termination rules and adaptive parameters.
- (c) As an upshot, we customize our algorithm to solve general constrained problems of the form (2). We prove the same $\mathcal{O}\left(\frac{1}{k}\right)$ -convergence rate guarantee on both the objective residual $|f(x^k) - f(x^*)|$ and the feasibility violation $\text{dist}_{\mathcal{K}}(Ax^k - b)$. This rate is given on the last iterate of the outer sequence.

Our algorithm is a primal-dual method, which can solve composite convex problem with linear operators as in Chambolle–Pock's method [13]. It only requires one proximal operator of f and g^* , one matrix-vector multiplication and one adjoint for each iteration. It is parallelizable when f is separable, i.e., $f(x) = \sum_{i=1}^N f_i(x_{[i]})$. Under this structure, our method has more advantages than ADMM and Chambolle–Pock's method. In the algorithm, we provide explicit rules to update all algorithmic parameters. We also note that these updates can be modified to trade-off between the primal or the dual progress. Since the parameters in ASGARD [59] are decreasing making its step-size smaller at each iteration, the restarting schemes developed in this paper reset these parameters to preserve large step-size at each outer-loop iteration.

Paper organization: The rest of this paper is organized as follows. Section 2 recalls some necessarily mathematical background and tools. Section 3 presents our main result with algorithm and its convergence guarantee. We study both unconstrained and constrained cases. Section 3.4 shows an extension of our method to three composite objective functions with linearization on potentially smooth terms. In Sect. 4, we provide seven numerical examples to test our algorithm against state-of-the-arts. Section 5 compares our method and existing algorithms in the literature.

2 Mathematical tools

We review some key ingredients for the design of our primal-dual methods.

Notation: We denote the norm in the primal space \mathcal{X} by $\|\cdot\|_{\mathcal{X}}$ and the norm in the dual space \mathcal{Y} by $\|\cdot\|_{\mathcal{Y}}$. Their dual norms are denoted by $\|\cdot\|_{\mathcal{X},*}$ and $\|\cdot\|_{\mathcal{Y},*}$, respectively. Given a positive real number a , $\lfloor a \rfloor$ denotes the largest integer less than or equal to a .

For a given nonempty, closed, and convex set \mathcal{K} in \mathbb{R}^p , we denote $\delta_{\mathcal{K}}(x) = 0$, if $x \in \mathcal{K}$, and $\delta_{\mathcal{K}}(x) = +\infty$, otherwise, its indicator function; and $s_{\mathcal{K}}(y) := \sup_{x \in \mathcal{K}} \langle x, y \rangle$ its support function. We define the normal cone of \mathcal{K} as $\mathcal{N}_{\mathcal{K}}(x) := \{w \in \mathbb{R}^n \mid \langle w, x - y \rangle \geq 0, y \in \mathcal{K}\}$ if $x \in \mathcal{K}$; $\mathcal{N}_{\mathcal{K}}(x) := \emptyset$, otherwise. We also define $\mathcal{K}^o := \{w \in \mathbb{R}^n \mid \langle w, x \rangle \leq 1, x \in \mathcal{K}\}$ as the polar set of \mathcal{K} . If \mathcal{K} is a convex cone, then $\mathcal{K}^o = -\mathcal{K}^*$, where $\mathcal{K}^* := \{w \in \mathbb{R}^n \mid \langle w, x \rangle \geq 0, x \in \mathcal{K}\}$ is the dual cone of \mathcal{K} .

Given a proper, closed, and convex function $f : \mathbb{R}^p \rightarrow \mathbb{R} \cup \{+\infty\}$, we use $\text{dom}(f) := \{x \in \mathbb{R}^p \mid f(x) < +\infty\}$ to denote its domain and $\partial f(x) := \{w \in \mathbb{R}^p \mid f(y) \geq f(x) + \langle w, y - x \rangle, y \in \text{dom}(f)\}$ to denote its subdifferential at x . If f is differentiable, $\nabla f(x)$ denotes its gradient at x . As usual, $f^*(y) := \sup_x \{\langle x, y \rangle - f(x)\}$ denotes the Fenchel conjugate of f . We say that $f : \mathcal{X} \rightarrow \mathbb{R}$ has L_f -Lipschitz gradient if it satisfies $\|\nabla f(x) - \nabla f(y)\|_{\mathcal{X},*} \leq L_f \|x - y\|_{\mathcal{X}}$, for any $x, y \in \mathcal{X}$. We say that f is Lipschitz continuous on \mathcal{X} with a Lipschitz constant $M_f \in [0, +\infty)$ if $|f(x) - f(y)| \leq M_f \|x - y\|_{\mathcal{X}}$ for any $x, y \in \mathcal{X}$. We use $\text{prox}_f(x) := \text{argmin}_u \{f(u) + (1/2)\|u - x\|_{\mathcal{X}}^2\}$ to denote the proximal operator of f . We say that f is “proximally tractable” if prox_f can be computed efficiently, e.g., in a closed form, or by a polynomial algorithm. By Moreau’s identity, we have $\text{prox}_{\gamma f}(x) + \gamma \text{prox}_{\gamma^{-1} f^*}(\gamma^{-1}x) = x$ for any $x \in \text{dom}(f)$.

2.1 Primal-dual formulation

Dual problem and min–max formulation: Associated with the primal problem (1), we also consider the corresponding dual problem:

$$D^* := \min_{y \in \mathbb{R}^n} \left\{ D(y) := f^*(-A^\top y) + g^*(y) \right\}, \tag{4}$$

where f^* and g^* are the Fenchel conjugates of f and g , respectively. Clearly, we can write the pair (1)–(4) into the following min–max saddle-point problem:

$$\min_{x \in \mathbb{R}^p} \max_{y \in \mathbb{R}^n} \left\{ \mathcal{L}(x, y) := f(x) + \langle Ax, y \rangle - g^*(y) \right\}. \tag{5}$$

Under mild and standard assumptions, this min–max problem is solvable and achieves zero duality gap, i.e., $P^* + D^* = 0$. In particular, the dual problem of (2) can be written as follows:

$$D^* := \min_{y \in \mathbb{R}^n} \left\{ D(y) := f^*(-A^\top y) + \langle b, y \rangle + s_{\mathcal{K}}(y) \right\}, \tag{6}$$

where $s_{\mathcal{K}}(y) = \sup_{x \in \mathcal{K}} \langle y, x \rangle$ is the support function of \mathcal{K} . Compared to (4), we have $g^*(y) = \langle b, y \rangle + s_{\mathcal{K}}(y) = s_{b+\mathcal{K}}(y)$. Let \mathcal{X}^* and \mathcal{Y}^* be the solution sets of the primal problem (1) [or (2)] and dual problem (4) [or (6)], respectively.

Fundamental assumptions: Throughout this paper, we will develop methods for solving (1) and (2). We propose a unified set of assumptions that covers both unconstrained and constrained problems in (1) and (2), respectively.

Assumption 1 *We impose the following assumptions on (1):*

1. The solution set \mathcal{X}^* of (1) is nonempty.
2. Both f and g are proper, closed, and convex.
3. The function g is Lipschitz continuous its domain: there exists $\widehat{M}_g \in (0, +\infty)$ such that $|g(x) - g(y)| \leq \widehat{M}_g \|x - y\|_{\mathcal{X}}$ for any $x, y \in \text{dom}(g)$.
4. The Slater condition $\text{ri}(\text{dom}(f)) \cap \{x \in \mathbb{R}^p \mid Ax \in \text{ri}(\text{dom}(g))\} \neq \emptyset$ holds, where $\text{ri}(\mathcal{X})$ is the relative interior of \mathcal{X} .

The details of Slater conditions can be found in, e.g., [2]. Assumption 1 covers the case of equality constraint, i.e., $g = \delta_{\{b\}}$, cone constraints, i.e. $g = \delta_{\mathcal{K}}$ as in (2), where \mathcal{K} is a convex cone, Lipschitz continuous functions, i.e. g is globally Lipschitz continuous, and combinations thereof. It guarantees the strong duality of (1) and (4) to hold.

Optimality conditions: Associated with the primal and dual problems (1)–(4), we have the following optimality conditions:

$$0 \in \partial f(x^*) + A^\top \partial g(Ax^*) \quad \text{and} \quad 0 \in -A \partial f^*(-A^\top y^*) + \partial g^*(y^*).$$

We can write this optimality condition into the following KKT condition:

$$0 \in \partial f(x^*) + A^\top y^* \quad \text{and} \quad 0 \in -Ax^* + \partial g^*(y^*). \tag{7}$$

For the constrained problem (2), these conditions are written as

$$0 \in \partial f(x^*) + A^\top y^*, \quad Ax^* - b \in \mathcal{K}, \quad \text{and} \quad y^* \in \mathcal{N}_{\mathcal{K}}(Ax^* - b),$$

where $\mathcal{N}_{\mathcal{K}}(\cdot)$ is the normal cone of \mathcal{K} defined above. If \mathcal{K} is a closed, pointed, and convex cone, then $\mathcal{N}_{\mathcal{K}} \equiv -\mathcal{K}^*$ the dual cone of \mathcal{K} . In this case, $y^* \in -\mathcal{K}^*$.

2.2 Bregman distances and generalized proximal operators

In the sequel, we will use Bregman distances for smoothing and computing proximal operators. Therefore, we give basic properties on Bregman distances.

Let $p_{\mathcal{Z}}$ be μ_p -strongly convex, continuous, and differentiable on \mathcal{Z} with the strong convexity $\mu_p = 1$, where $\mathcal{Z} = \text{dom}(p_{\mathcal{Z}}) \neq \emptyset$. We call $p_{\mathcal{Z}}$ a proximity function (or prox-function). We define the Bregman distance induced by $p_{\mathcal{Z}}$ as

$$b_{\mathcal{Z}}(x, y) := p_{\mathcal{Z}}(x) - p_{\mathcal{Z}}(y) - \langle \nabla p_{\mathcal{Z}}(y), x - y \rangle, \quad \forall x, y \in \mathcal{Z}.$$

As special cases, if we choose $p_{\mathcal{Z}}(x) = \frac{1}{2} \|x\|_2^2$, then $b_{\mathcal{Z}}(x, y) = \frac{1}{2} \|x - y\|_2^2$, the standard Euclidean distance square. If we choose $p_{\mathcal{Z}}(x) := \sum_i x_i \ln(x_i)$, the entropy function for a standard simplex Δ , then $b_{\mathcal{Z}}(x, y) := \sum_i x_i \ln\left(\frac{x_i}{y_i}\right) - x_i + y_i$, the so-called Kullback–Leibler (KL) divergence. Moreover, it is obvious that $b_{\mathcal{Z}}(x, y) \geq \frac{1}{2} \|x - y\|_{\mathcal{Z}}^2$ for all $x, y \in \mathcal{Z}$. When a Bregman distance $b_{\mathcal{Z}}$ has Lipschitz continuous gradient, we denote its Lipschitz constant by $L_{b_{\mathcal{Z}}}$. We refer to [16,25,36] for several concrete examples of Bregman divergences.

2.3 Nesterov’s smoothing technique

We focus on Nesterov’s smoothing technique with general Bregman distances [4,48] as follows. Since g in (1) is possibly nonsmooth, we smooth it by

$$g_\beta(u; \dot{y}) := \max_{y \in \mathcal{Y}} \left\{ \langle u, y \rangle - g^*(y) - \beta b_{\mathcal{Y}}(y, \dot{y}) \right\}, \tag{8}$$

where $\dot{y} \in \mathbb{R}^n$ is the prox-center point, and $\beta > 0$ is a smoothness parameter. The function $g_\beta(\cdot; \dot{y})$ is convex and smooth, its gradient is given by

$$\nabla g_\beta(u; \dot{y}) = y_\beta^*(u; \dot{y}) = \operatorname{argmin}_{y \in \mathcal{Y}} \left\{ g^*(y) - \langle u, y \rangle + \beta b_{\mathcal{Y}}(y, \dot{y}) \right\}. \tag{9}$$

Clearly, $\nabla g_\beta(\cdot; \dot{y})$ is Lipschitz continuous with the Lipschitz constant $L_{g_\beta} = \frac{1}{\beta}$. Moreover, we have

$$g_\beta(u; \dot{y}) \leq g(u) \leq g_\beta(u; \dot{y}) + \beta D_{\mathcal{Y}}, \tag{10}$$

where $D_{\mathcal{Y}} := \sup \{ b_{\mathcal{Y}}(y, \dot{y}) \mid y \in \operatorname{dom}(g^*) \}$ is the prox-diameter of g^* . Here, $D_{\mathcal{Y}}$ is finite if and only if g is Lipschitz continuous on $\operatorname{dom}(f)$ with the Lipschitz constant $L_g := \sqrt{2}D_{\mathcal{Y}}$, i.e., $|g(u) - g(v)| \leq \sqrt{2}D_{\mathcal{Y}}\|u - v\|$ for all $u, v \in \operatorname{dom}(g)$ due to [7, Proposition 4.4.6].

If we choose $b_{\mathcal{Y}}(y, \dot{y}) = \frac{1}{2}\|y - \dot{y}\|_2^2$, then we can write $y_\beta^*(u; \dot{y})$ as:

$$\nabla g_\beta(u; \dot{y}) = \operatorname{argmin}_{y \in \mathbb{R}^n} \left\{ g^*(y) - \langle u, y \rangle + \frac{\beta}{2}\|y - \dot{y}\|^2 \right\} = \operatorname{prox}_{g^*/\beta} \left(\dot{y} + \frac{1}{\beta}u \right). \tag{11}$$

Smoothing techniques are widely used in the literature, including [4,8,9,23,44]. The idea of smoothing is to approximate the original problem (1) by a (partially) smoothed problem. For example, in our setting, we smooth g and consider the following smoothed problem:

$$P_\beta^* := \min_{x \in \mathbb{R}^p} \left\{ P_\beta(x; \dot{y}) := f(x) + g_\beta(Ax; \dot{y}) \right\}. \tag{12}$$

We define the following generalized proximal operator with Bregman distance $d_{\mathcal{X}}$ induced by a prox-function $q_{\mathcal{X}}$:

$$\mathcal{P}_{\theta f}^{d_{\mathcal{X}}}(u, y) := \operatorname{argmin}_{v \in \mathcal{X}} \left\{ f(v) + \langle y, v - u \rangle + \frac{1}{\theta}d_{\mathcal{X}}(v, u) \right\}. \tag{13}$$

Given that the Bregman distance $d_{\mathcal{X}}$ is defined in \mathcal{X} and $b_{\mathcal{Y}}$ is defined in \mathcal{Y} , we define the following operator norm of A :

$$\|A\| := \max_{x \in \mathbb{R}^p} \left\{ \frac{\|Ax\|_{\mathcal{Y},*}}{\|x\|_{\mathcal{X}}} \right\}. \tag{14}$$

Different from [4,9,23,44,48], our strategy allows one to update the smoothness parameter β gradually at each iteration. Similar work can be found in [8,47], which are also essentially different from ours as discussed in Sect. 5.

3 Main results: self-adaptive double-loop ASGARD

In this section, we develop a self-adaptive double-loop accelerated smoothed primal-dual gap reduction algorithm to solve (1) and (2). We first present the complete algorithm. Next, we provide its convergence analysis. Then, we specify our algorithm to handle the constrained setting (2). Finally, we extend our method to handle (1) with the sum of three objective functions where the third function has Lipschitz gradient.

3.1 The algorithm and its convergence guarantee

Main idea: The proposed algorithm consists of two loops:

- The inner loop performs an accelerated proximal gradient (APG) scheme [60] to solve the smoothed problem (12) for a fixed β , which is different from [59], where β is updated at each iteration. We note that in the constrained case, the smoothed problem (12) is the augmented Lagrangian formulation.
- The outer loop can be considered as a restarting step and simultaneously decreases the smoothness parameter β .

The intuition behind our new strategy lies on the fact that when applied to (12) with a fixed β , APG gets $\mathcal{O}\left(\frac{1}{k^2}\right)$ rate, whereas ASGARD as presented in [59] controls the parameters in such a way that the algorithm gets $\mathcal{O}\left(\frac{1}{k}\right)$ rate throughout its execution. The idea is to take the advantage of the faster rate of APG for the inner loop while carefully adjusting the number of inner iterations and the smoothness parameter to get the same overall $\mathcal{O}\left(\frac{1}{k}\right)$ rate with better practical performance. Our analysis also gives insights on the heuristic restart strategy outlined in [59]. For the sake of presentation and its flexibility of using Bregman distances in proximal operators, we choose Tseng's variant of APG [60]. However, we can replace it by another scheme such as FISTA [3]. We adaptively determine the number of inner iterations at each outer iteration. Therefore, there is no need to tune this parameter. The outer loop gradually decreases the smoothness parameter β such that the algorithm is still guaranteed to converge to the true solution of (1) or (2).

The algorithm: The complete algorithm is presented in Algorithm 1.

Algorithm 1 uses APG with **Option 1** at Step 8. This step requires one subproblem in \tilde{y}^{k+1} , one prox_f of f , one matrix-vector multiplication Ax and its adjoint $A^\top y$ at each iteration. Hence, Algorithm 1 has the same per-iteration complexity as ASGARD, except for the extra step, Step 14, where we update the dual center \tilde{y}^s at each outer loop iteration. In general, the number of outer iterations is small as it is the number of restarting steps. Hence, Step 14 does not significantly increase the

Algorithm 1 (Self-Adaptive Double Loop ASGAR Algorithm)

- 1: **Initialization:**
 - 2: Choose $\beta_0 > 0$, $\omega > 1$, a positive integer $m_0 \geq 1$, $\bar{x}^0 \in \mathbb{R}^p$, and $\dot{y}^0 \in \mathbb{R}^n$.
 - 3: Choose a Bregman distance $b_{\mathcal{Y}}$ for y and $d_{\mathcal{X}}$ for x .
 - 4: Set $K_0 \leftarrow 0$, $\hat{x}^0 \leftarrow \bar{x}^0$, and $\tau_0 \leftarrow 1$.
 - 5: **For** $s := 0$ **to** $S_{\max} - 1$, **perform:**
 - 6: **For** $j := 0$ **to** $m_s - 1$ **perform**
 - 7: Set $k \leftarrow K_s + j$.
 - 8: Update
$$\begin{cases} \bar{x}^k \leftarrow (1 - \tau_k)\bar{x}^k + \tau_k\hat{x}^k \\ \tilde{y}^{k+1} \leftarrow \operatorname{argmin}_{y \in \mathcal{Y}} \left\{ g^*(y) - \langle A\bar{x}^k, y \rangle + \beta_s b_{\mathcal{Y}}(y, \dot{y}^s) \right\} \\ \hat{x}^{k+1} \leftarrow \mathcal{P}_{\gamma_k f}^{d_{\mathcal{X}}}(\hat{x}^k, A^\top \tilde{y}^{k+1}) \quad \text{with } \gamma_k \leftarrow \frac{\beta_s}{\|A\|^2 \tau_k}. \end{cases}$$
 - 9: Update \bar{x}^{k+1} using one of the following two options:

$$\begin{cases} \bar{x}^{k+1} \leftarrow \bar{x}^k + \tau_k(\hat{x}^{k+1} - \hat{x}^k) & \text{(Option 1: Averaging step)} \\ \bar{x}^{k+1} \leftarrow \mathcal{P}_{\beta_s f / \|A\|^2}^{d_{\mathcal{X}}}(\bar{x}^k, A^\top \tilde{y}^{k+1}) & \text{(Option 2: Proximal step).} \end{cases}$$
 - 10: Update $\tau_k \leftarrow \frac{2}{k - K_s + 2}$.
 - 11: **End For**
 - 12: Update $K_{s+1} \leftarrow K_s + m_s$.
 - 13: Restart $\bar{x}^{K_{s+1}} \leftarrow \hat{x}^{K_{s+1}} \equiv \hat{x}^{K_s + m_s}$.
 - 14: Restart $\dot{y}^{s+1} \leftarrow \operatorname{prox}_{\frac{1}{\beta_s} g^*} \left(\dot{y}^s + \frac{1}{\beta_s} A\bar{x}^{K_{s+1}} \right)$.
 - 15: Restart $\tau_{K_{s+1}} \leftarrow 1$.
 - 16: Set $m_{s+1} \leftarrow \lfloor \omega(m_s + 1) + 1 \rfloor - 1$.
 - 17: Set $\beta_{s+1} \leftarrow \frac{\beta_s(m_{s+1} + 1)}{\omega \sqrt{m_{s+1}(m_{s+1} + 3)}}$
 - 18: **End For**
-

overall computational cost of the entire algorithm. Note that \bar{x}^{k+1} computed at Step 9 using **Option 1** is a weighted averaging step. To avoid this averaging, we can choose **Option 2**, which requires an additional generalized proximal operator of f .

We can replace Step 8 of Algorithm 1 by the following FISTA step:

$$\begin{cases} \tilde{y}^{k+1} \leftarrow \operatorname{argmin}_{y \in \mathcal{Y}} \left\{ g^*(y) - \langle A\bar{x}^k, y \rangle + \beta_s b_{\mathcal{Y}}(y, \dot{y}^s) \right\} \\ \bar{x}^{k+1} \leftarrow \mathcal{P}_{\gamma_k f}^{d_{\mathcal{X}}}(\bar{x}^k, A^\top \tilde{y}^{k+1}) \quad \text{with } \gamma_k \leftarrow \frac{\beta_s}{\|A\|^2} \\ \tilde{x}^{k+1} \leftarrow \bar{x}^{k+1} + \frac{(1 - \tau_k)\tau_{k+1}}{\tau_k}(\bar{x}^{k+1} - \bar{x}^k). \end{cases} \tag{15}$$

However, we need to replace the general Bregman distance $d_{\mathcal{X}}$ by an Euclidean distance $d_{\mathcal{X}}(\cdot, \dot{x}) := \frac{1}{2} \|\cdot - \dot{x}\|_2^2$. The scheme (15) allows us to compute \bar{x}^k through $\mathcal{P}_{\gamma_k f}^{d_{\mathcal{X}}}(\cdot)$ instead of a weighted averaging step as with **Option 1**.

Comparison between ASGARD [59] and Algorithm 1: When designing Algorithm 1, we focused on improving the practical efficiency of ASGARD while retaining the $\mathcal{O}(\frac{1}{k})$ -worst-case rate on the last primal iterate for the objective residual and feasibility violation. To achieve this goal, we introduced several fundamental changes to the original ASGARD algorithm.

Firstly, ASGARD only works with the Euclidean distance in the primal space while Algorithm 1 works with any Bregman distance.

Secondly, accelerated proximal gradient or (15) serves as a sub-routine from Step 6 to Step 11 in Algorithm 1 when β_k is fixed at β_s , while (15) is intertwined with updates of β_k in ASGARD.

Thirdly, as we discussed earlier, the parameter τ_k in ASGARD is gradually decreased to zero, making its performance slow. Algorithm 1 allows one to reset τ_k back to one at Step 15 making use of larger step-sizes at Steps 8 and 9. We can view Algorithm 1 as applying a multiple stages strategy to ASGARD, where we rerun ASGARD at a new but better initial point at each stage s .

Finally, Algorithm 1 can serve as a unified framework for convergence analysis, where we can replace the inner loop with any other accelerated schemes such as stochastic and coordinate descent variants.

The following lemma provides a key estimate for the optimality condition of (2), whose proof is given in Appendix 6.2.

Lemma 1 *Suppose that $b_{\mathcal{Y}}$ has an $L_{b_{\mathcal{Y}}}$ -Lipschitz gradient, $L_{b_{\mathcal{Y}}} \in (0, +\infty]$ and that g is \widehat{M}_g -Lipschitz continuous on its domain. Let (x^*, y^*) be a saddle-point of the Lagrange function of (1) and*

$$S_{\beta}(\bar{x}, \dot{y}) := \max_{y \in \mathbb{R}^n} \left\{ f(\bar{x}) + \langle y, A\bar{x} \rangle - g^*(y) - \beta b_{\mathcal{Y}}(y, \dot{y}) - f(x^*) - g(Ax^*) \right\}. \tag{16}$$

Let $\beta_b := \beta L_{b_{\mathcal{Y}}}$, $\bar{y}_{\beta}^ = y_{\beta}^*(A\bar{x}, \dot{y})$, and \bar{z} be the projection of $A\bar{x}$ onto $\text{dom}(g)$. If either $L_{b_{\mathcal{Y}}} < +\infty$ or $D_{\mathcal{Y}} < +\infty$, then we have*

$$\left\{ \begin{array}{l} f(\bar{x}) + g(\bar{z}) - P^* \geq -\|y^*\|_{\mathcal{Y}} \text{dist}_{\mathcal{Y},*}(A\bar{x}, \text{dom}(g)) \\ f(\bar{x}) + g(\bar{z}) - P^* \leq S_{\beta}(\bar{x}, \dot{y}) \\ \qquad \qquad \qquad + \beta \min \left\{ D_{\mathcal{Y}}, (2\widehat{M}_g + \|\dot{y}\|_{\mathcal{Y}}) \|\nabla_{\mathcal{Y}} b_{\mathcal{Y}}(\bar{y}_{\beta}^*, \dot{y})\|_{\mathcal{Y},*} \right\} \\ \text{dist}_{\mathcal{Y},*}(A\bar{x}, \text{dom}(g)) \leq \beta \|\nabla_{\mathcal{Y}} b_{\mathcal{Y}}(\bar{y}_{\beta}^*, \dot{y})\|_{\mathcal{Y},*} \\ \qquad \qquad \qquad \leq \beta_b \left[\|y^* - \dot{y}\|_{\mathcal{Y}} + \sqrt{\|y^* - \dot{y}\|_{\mathcal{Y}}^2 + \frac{2}{\beta_b} S_{\beta}(\bar{x}, \dot{y})} \right], \end{array} \right. \tag{17}$$

where $\text{dist}_{\mathcal{Y},*}(u, \mathcal{C}) := \inf_{v \in \mathcal{C}} \|u - v\|_{\mathcal{Y},*}$ is the distance from u to \mathcal{C} .

Now, we are ready to state the main convergence result in the following theorem, whose proof is given in Appendix 7.

Theorem 1 Assume that Assumption 1 holds. Let $\{\bar{x}^{K_s}\}$ be the sequence generated by Algorithm 1 and \bar{z}^{K_s} be the projection of $A\bar{x}^{K_s}$ onto $\text{dom}(g)$. If either $L_{b_Y} < +\infty$ or $D_Y < +\infty$, then we have

$$\begin{cases} f(\bar{x}^{K_s}) + g(\bar{z}^{K_s}) - P^* \geq -\|y^*\|_Y \text{dist}_{Y,*}(A\bar{x}^{K_s}, \text{dom}(g)) \\ f(\bar{x}^{K_s}) + g(\bar{z}^{K_s}) - P^* \leq \frac{\omega\kappa_0 R_0^2}{\rho_0 [(\omega - 1)K_s + \kappa_0]} + \frac{\beta_0\omega\kappa_0 C_0^*}{(\omega - 1)K_s + \kappa_0} \\ \text{dist}_{Y,*}(A\bar{x}^{K_s}, \text{dom}(g)) \leq \frac{\beta_0 L_{b_Y} \omega \kappa_0 (\sqrt{2} + 2) R_0}{\rho_0 [(\omega - 1)K_s + \kappa_0]}, \end{cases} \quad (18)$$

where y^* is any dual solution of (4), and

$$\begin{cases} \rho_0 := \beta_0 \left(1 - \frac{1}{(\omega-1)m_0}\right) \\ \kappa_0 := m_0 + \frac{\omega}{\omega-1} \\ R_0 := \left[\frac{4\|A\|^2}{(m_0+1)^2} d_{\mathcal{X}}(x^*, \bar{x}^0) + \frac{\beta_0^2 m_0(m_0+3)}{(m_0+1)^2} b_Y(y^*, y^0) \right]^{1/2} \\ C_0^* := \min \left\{ D_Y, (\sqrt{2} + 2) \frac{L_{b_Y} R_0}{\rho_0} [2\widehat{M}_g + \|y^*\|_Y + \frac{\sqrt{2}R_0}{\rho_0}] \right\}. \end{cases}$$

Consequently, Algorithm 1 achieves an $\mathcal{O}(\frac{1}{K_s})$ convergence rate in a **non-ergodic sense**, i.e. the objective satisfies $|f(\bar{x}^{K_s}) + g(\bar{z}^{K_s}) - P^*| \leq \mathcal{O}(\frac{1}{K_s})$ and the constraints satisfy $\text{dist}_{Y,*}(A\bar{x}^{K_s}, \text{dom}(g)) \leq \mathcal{O}(\frac{1}{K_s})$.

Remark 1 We make a few remarks about Theorem 1.

1. The smoothness parameter β is only updated at the outer loop but with a geometric rate, depending on the factor parameter ω . We can select different ω to observe its performance in particular applications.
2. The convergence rate is given at the last iterate instead of the averaged sequence as often seen in other primal-dual methods [13,51,56,66].
3. In the case where g is Lipschitz continuous ($D_Y < +\infty$), we can choose any dual prox-function. Indeed, as $\text{dist}_{Y,*}(A\bar{x}^{K_s}, \text{dom}(g)) = 0$, we do not need the third inequality and we can choose b_Y such that $L_{b_Y} = +\infty$. Moreover, $\bar{z}^{K_s} = A\bar{x}^{K_s}$. Hence, $f(\bar{x}^{K_s}) + g(\bar{z}^{K_s}) = f(\bar{x}^{K_s}) + g(A\bar{x}^{K_s}) = P(\bar{x}^{K_s})$. See Corollary 1 below.
4. When \mathcal{Y} is unbounded, i.e. $D_Y = +\infty$ (for instance when there are constraints), we need to choose a smooth prox-function for the dual problem. We present this case in detail in Sect. 3.3.
5. Evaluating the projection of a vector onto $\text{dom}(g)$ is usually a simple task when prox_g has an explicit form.
6. The bound of convergence rate depends on both the prox-distance between \bar{x}^0 to x^* and y^0 to y^* .

Remark 2 By using (44) in our analysis, we can show that the objective residual $\{f(\bar{x}^{K_s}) + g(\bar{z}^{K_s})\}$ and the constraint violation $\text{dist}_{Y,*}(A\bar{x}^{K_s}, \text{dom}(g))$ sequences converge to P^* and zero, respectively at the rate of $\mathcal{O}(\frac{1}{k})$ for any $k \geq 1$ instead of

$k = K_s$ at the outer loop only. If we use the averaging step of APG, then \bar{x}^k is computed via a weighted averaging step of the inner loop. However, if we use the proximal step in APG or the FISTA scheme (15), then \bar{x}^k is computed through the generalized proximal operator $\mathcal{P}_{\beta_s f/\|A\|^2}^{d\mathcal{X}}$. This rate is fully non-ergodic for both inner and outer loops.

3.2 Application to Lipschitz convex optimization

If g is globally Lipschitz continuous on \mathbb{R}^n , then we recover the framework of [48]. In this case, $\text{dom}(g) = \mathbb{R}^n$ and $\text{dist}_{\mathcal{Y},*}(A\bar{x}^{K_s}, \text{dom}(g)) = 0$. Moreover, $D_{\mathcal{Y}} < +\infty$. Theorem 1 can be simplified as follows.

Corollary 1 *Assume that Assumption 1 holds and g is M_g -globally Lipschitz continuous. Let $\{\bar{x}^{K_s}\}$ be the sequence generated by Algorithm 1. Then*

$$0 \leq P(\bar{x}^{K_s}) - P^* \leq \frac{\omega\kappa_0}{(\omega - 1)K_s + \kappa_0} \left(\frac{R_0^2}{\rho_0} + \beta_0 D_{\mathcal{Y}} \right), \tag{19}$$

where ρ_0, κ_0 , and R_0 are defined as in Theorem 1.

3.3 Application to constrained convex optimization

In this subsection, we specify Algorithm 1 to solve the constrained problem (2). First, we choose the Bregman distance used for smoothing of the dual variables to have Lipschitz gradient. Under this condition, and note that \dot{y} is the prox-center of $b_{\mathcal{Y}}$, we have

$$b_{\mathcal{Y}}(y, \dot{y}) \leq \frac{Lb_{\mathcal{Y}}}{2} \|y - \dot{y}\|_{\mathcal{Y}}^2. \tag{20}$$

Let us define $g(Ax) := \delta_{\mathcal{K}}(Ax - b)$ the indicator function of \mathcal{K} , where the set \mathcal{K} is such that Slater’s condition in Assumption 1 holds. Then, we can write

$$g(Ax) := \sup_{y \in \mathbb{R}^n} \{ \langle Ax - b, y \rangle - s_{\mathcal{K}}(y) \}, \tag{21}$$

where $s_{\mathcal{K}}(y) := \sup_{u \in \mathcal{K}} \langle y, u \rangle$ is the support function of \mathcal{K} . In this case, the smoothed function $g_{\beta}(Ax; \dot{y})$ becomes

$$g_{\beta}(Ax; \dot{y}) := \max_{y \in \mathbb{R}^n} \left\{ \langle Ax - b, y \rangle - s_{\mathcal{K}}(y) - \beta b_{\mathcal{Y}}(y, \dot{y}) \right\}. \tag{22}$$

Example 1 Assume that we choose $b_{\mathcal{Y}}(x, \dot{x}) = \frac{1}{2} \|x - \dot{x}\|_2^2$. Then

$$g_{\beta}(Ax; \dot{y}) = \frac{1}{2\beta} \text{dist}_{\mathcal{K}}(Ax - b + \beta\dot{y})^2 - \frac{\beta}{2} \|\dot{y}\|_2^2. \tag{23}$$

Moreover, the solution $y_\beta^*(Ax; \dot{y})$ of the maximization problem in (22) is

$$y_\beta^*(Ax; \dot{y}) = \dot{y} + \frac{1}{\beta} (Ax - b - \text{proj}_{\mathcal{K}} (Ax - b + \beta \dot{y})), \tag{24}$$

where $\text{proj}_{\mathcal{K}}(\cdot)$ denotes the projection onto \mathcal{K} .

In particular, if \mathcal{K} is a cone, then $y_\beta^*(Ax; \dot{y}) = \text{proj}_{-\mathcal{K}^*} \left(\dot{y} + \frac{1}{\beta} (Ax - b) \right)$, where \mathcal{K}^* is the dual cone of \mathcal{K} . The dual step for computing \tilde{y}^k at the second line of Step 8 of Algorithm 1 becomes

$$\begin{aligned} \tilde{y}^{k+1} &\leftarrow \dot{y}^s + \frac{1}{\beta_s} \left(A\tilde{x}^k - b - \text{proj}_{\mathcal{K}} \left(A\tilde{x}^k - b + \beta_s \dot{y}^s \right) \right) \\ &= \frac{1}{\beta_s} \text{proj}_{-\mathcal{K}^*} \left(A\tilde{x}^k - b + \beta_s \dot{y}^s \right). \end{aligned} \tag{25}$$

In this case, we can apply Theorem 1 to (2) with $g(\bar{z}^{K_s+1}) = 0$, $D_{\mathcal{Y}} = +\infty$, and $\widehat{M}_g = 0$ to obtain the following corollary.

Corollary 2 *Assume that Assumption 1 holds and that $g(z) := \delta_{\mathcal{K}}(z - b)$. Let us choose $b_{\mathcal{Y}}$ such that (20) holds with $L_{b_{\mathcal{Y}}} < +\infty$. Let $\{\bar{x}^{K_s}\}$ be the sequence generated by Algorithm 1. Then, we have*

$$\begin{cases} f(\bar{x}^{K_s}) - P^* \geq -\|y^*\|_{\mathcal{Y}} \text{dist}_{\mathcal{Y},*} (A\bar{x}^{K_s}, \mathcal{K}) \\ f(\bar{x}^{K_s}) - P^* \leq \frac{\omega\kappa_0}{(\omega - 1)K_s + \kappa_0} \left[\frac{R_0^2}{\rho_0} + \frac{(\sqrt{2} + 2)\beta_0 L_{b_{\mathcal{Y}}} R_0}{\rho_0} \left(\|y^*\|_{\mathcal{Y}} + \frac{\sqrt{2}R_0}{\rho_0} \right) \right] \\ \text{dist}_{\mathcal{Y},*} (A\bar{x}^{K_s}, \mathcal{K}) \leq \frac{\beta_0 L_{b_{\mathcal{Y}}} \omega \kappa_0 (\sqrt{2} + 2) R_0}{\rho_0 [(\omega - 1)K_s + \kappa_0]}, \end{cases}$$

where y^* is any dual solution of (6), and ρ_0 , κ_0 , and R_0 are defined as in Theorem 1.

3.4 Extension to composite case with three objective terms

It is straightforward to apply Algorithm 1 in the presence of a smooth term in the objective. The problem template we focus on in this section is

$$F^* := \min_{x \in \mathbb{R}^p} \left\{ F(x) := f(x) + g(Ax) + h(x) \right\}, \tag{26}$$

where f and g are as described in Assumption 1 and h is a differentiable function with L_h -Lipschitz gradient. In this case, only Step 8 in Algorithm 1 needs to be modified as follows (see also in [49]):

$$\hat{x}^{k+1} \leftarrow \mathcal{P}_{\gamma_k f}^{d\mathcal{X}} \left(\hat{x}^k, \nabla h(\hat{x}^k) + A^\top \hat{y}^{k+1} \right) \quad \text{with } \gamma_k \leftarrow \frac{\beta_s}{\tau_k (\|A\|^2 + \beta_s L_h)}.$$

Note that this modification only changes the analysis of the inner loop as in [49] which does not affect our analysis of the outer loop. In addition, using L_h in the stepsize is not restrictive. When the Lipschitz constant is not known, line search strategies can be employed, see [49] for more details. The convergence of this variant is still guaranteed by Theorem 1 but the quantity R_0^2 will depend on L_h . We omit the details of this result here for succinctness.

4 Numerical experiments

We will test standard ASGARD [49,59], ASGARD with restart [49,59], and standard Chambolle–Pock’s algorithm [13] on the following problems. Note that when there is a smooth term in the objective, we use the version of Chambolle–Pock which linearizes the smooth term. This is also known in the literature as Vu–Condat’s algorithm [20,61]. We only compare with HOPS [68] in the first example because it does not apply to Basis pursuit, sparse subspace clustering, linear programming, and Markowitz’s portfolio optimization problems due to the unboundedness of the dual domain. For the ℓ_1 -SVM example, we observed that it is extremely slow and difficult to tune hyperparameters for different datasets. In all the experiments, we have used the standard distances $b_{\mathcal{Y}}(y_1, y_2) = \frac{1}{2} \|y_1 - y_2\|_2^2$ and $d_{\mathcal{X}}(x_1, x_2) = \frac{1}{2} \|x_1 - x_2\|_2^2$ for smoothing and computing the proximal operators for fair comparison with other methods which do not allow Bregman distances. In the sequel, we refer to our algorithm as ASGARD-DL, and we only run Option 1 at Step 9. In some cases, we also compare with ADMM and its variants.

The parameters are set as follows. For Chambolle–Pock’s method, we set its stepsizes $\sigma := \frac{1}{\|A\|}$ and $\tau := \frac{0.9999}{\|A\|^2 \sigma}$, where A is the linear operator in (1). For ASGARD-DL, we choose $\omega := 1.2$ and $m_0 := 6$, which give us comparable performance. For restarting ASGARD, we set the restarting frequency to be $s = 10$ in all experiments.

4.1 Convergence guarantees: ergodic versus non-ergodic rates

ALM, ADMM, and Chambolle–Pock methods have convergence rate guarantees in ergodic sense. That is, they have such a rate guarantee only on the averaged sequence. In contrast, our guarantee on the last iterate sequence generated by the algorithm. To illustrate the importance between these two, we consider two synthetic problems in this section. The first one is a square-root LASSO problem widely studied in the literature, which is given by:

$$F^* := \min_{x \in \mathbb{R}^p} \left\{ F(x) := \|Ax - b\|_2 + \lambda \|x\|_1 \right\}, \tag{27}$$

where $A \in \mathbb{R}^{n \times p}$ is generated using a Gaussian distribution and is normalized such that column norms are equal to 1. Given a ground-truth vector x^\natural , we generate observations

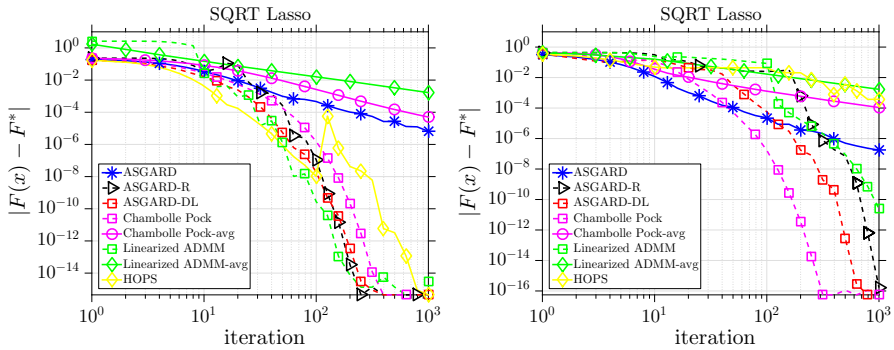


Fig. 1 The performance of 5 algorithms for solving the square-root LASSO problem (27). Left: $\sigma = 0.1, \lambda = 0.04$, right: $\sigma = 0.01, \lambda = 0.03$

as $b := Ax^\natural + \sigma \mathbf{n}$, where \mathbf{n} is a noise vector generated by a standard Gaussian distribution and $\sigma = 0.01$. We set $\lambda = 0.03$, which is tuned to get a good recovery of x^\natural .

In this experiment, we test the ergodic and non-ergodic variants of Linearized ADMM (in the sense that the augmented term in the Lagrangian is linearized) [30] and Chambolle–Pock’s algorithm [13], as well as the primal-dual homotopy smoothing method, called HOPS [68]. The methods in [13,30] have convergence guarantees for their last iterates, however, their rate guarantees only apply to the averaged sequence. Moreover, they are very successful to solve this type of problems as can be seen from the literature. The behavior of the algorithms is given in Fig. 1.

As can be seen from Fig. 1, the last iterate sequences of the Linearized ADMM and Chambolle Pock’s algorithms seem to have the best performance. However, the averaged iterate sequence for which the methods have the best-known convergence rate guarantee shows the slowest convergence behavior. Our method has the same rate as restarted ASGARD, called ASGARD-R. Note that ASGARD-R does not have any convergence guarantees.

As can be observed from Fig. 1, the performance of HOPS shows too much fluctuation with different datasets. We suspect that the reason making HOPS fluctuated is that it requires 3 separate parameters which are very difficult to tune. One choice of parameters may work well for one dataset but it can perform very poorly for another. Note that we are using the same set of parameters for different datasets in all the algorithms. HOPS requires knowing $\epsilon_0 \geq F(x^0) - F^*$ which we bound it using the fact that $F^* \geq 0$ and we set $\epsilon_0 = F(x^0)$. The second parameter is the rate at which they decrease the smoothness parameter, which is similar to the ω parameter in our algorithm. The last parameter is the number of inner iterations. In contrast, Algorithm 1 only requires setting the initial parameters for the inner iteration and smoothness parameter. Then, it automatically updates them during the iterative process. As we illustrate, we have obtained a similar performance across different datasets with the same set of parameters for our method.

To illustrate the behavior of the last iterates of Linearized ADMM and Chambolle Pock’s algorithm, we consider a degenerate linear program which is also studied in [59]:

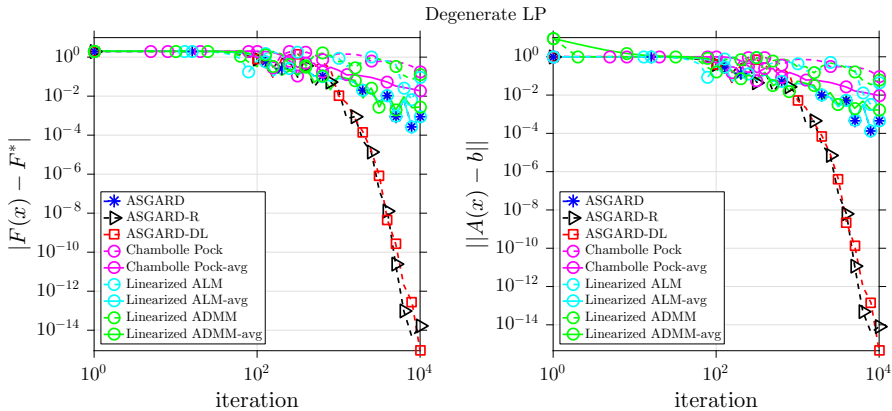


Fig. 2 The performance of 6 algorithms for solving the degenerate linear program (28)

$$\min_{x \in \mathbb{R}^p} \left\{ h(x) := 2x_p \mid \sum_{k=1}^{p-1} x_k = 1, x_p - \sum_{k=1}^{p-1} x_k = 0 \ (2 \leq j \leq n), x_p \geq 0 \right\}. \quad (28)$$

The second inequality is repeated $n - 1$ times which causes the problem to be degenerate. We define the linear constraint as

$$Ax := \left[\sum_{k=1}^{p-1} x_k, x_p - \sum_{k=1}^{p-1} x_k, \dots, x_p - \sum_{k=1}^{p-1} x_k \right]^T.$$

We have $b := (1, 0, \dots, 0)^T \in \mathbb{R}^n$. We map the problem to our template in (26) as $f(x) := \delta_{\{x_p \geq 0\}}(x_p)$, $g(x) := \delta_{\{b\}}(Ax)$, and $h(x) := 2x_p$. In this test, we choose $p = 10$ and $n = 200$.

In addition to Linearized ADMM and Chambolle–Pock’s algorithm, we also include linearized ALM [30] to solve this example. The result of this test is given in Fig. 2, where $F(x) = h(x)$.

As can be seen from Fig. 2, Linearized ADMM, Linearized ALM, and Chambolle–Pock’s algorithm can be extremely slow while our algorithm and ASGARD-R still make progress and converge to the optimal value with a very high accuracy, and beyond the theoretical worst-case rate guarantee.

4.2 Basis pursuit for recovering bag-of-words of text documents

We first consider a basis pursuit problem which is used in signal and image processing, statistics, and machine learning [12,18,24]:

$$\min_{x \in \mathbb{R}^p} \{ F(x) := \|x\|_1 \mid Ax = b \}, \quad (29)$$

where $A \in \mathbb{R}^{n \times p}$ and $b \in \mathbb{R}^n$. This problem clearly fits into our template (1) by mapping $f(\cdot) = \|\cdot\|_1$ and $g(\cdot) = \delta_{\{b\}}(\cdot)$. It is also a special case of (2) with $\mathcal{K} = \{\mathbf{0}\}$. The proximal operators of both terms are given in a closed form.

We apply this model to a text processing problem. In [1], the authors proposed using the basis pursuit formulation (29) to obtain bag-of-words representation from the unigram embedding representation of text. The setting can be briefly described as the following: For any word w , there exists a word vector $v_w \in \mathbb{R}^n$. For a given text document $\{w_1, \dots, w_T\}$, one defines the unigram embedding as $\sum_{i=1}^T v_{w_i}$. It is easy to see that the unigram embeddings can be written as a linear system Ax where $A \in \mathbb{R}^{n \times p}$ contains v_{w_i} in the i th column and $x \in \mathbb{R}^p$ is the bag-of-words vector which counts the number of occurrences of words in a text document. This application is considered in text processing applications to obtain the original text document given its unigram embeddings [63].

For this experiment, we have used the movie review dataset from [40]. We have selected 4 different documents and computed the unigram embeddings using pre-trained word embeddings from GloVe [53] with $n = 50$ as the dimension of the word vectors and restricted the vocabulary size to $p = 10,000$ for getting faster results with all algorithms.

We have applied 4 methods to solve (29) for 4 different documents. Here, the parameter β_0 in ASGARD, ASGARD-restart, and ASGARD-DL is set to $\beta_0 := 10\|A\|$. Note that this choice is not optimal, but give us reasonable results in all test. The results are compiled in Fig. 3.

As we can observe from Fig. 3, our new algorithm works quite well and is comparable with state-of-the-art methods for low accuracy. It outperforms them if we run the algorithms long enough to get a more accurate solution than $\varepsilon := 10^{-5}$ both in the objective residual and the feasibility violation. Note that (29) is fully nonsmooth, and A is non-orthogonal. If we apply ADMM to solve (29), then it requires to solve a general convex subproblem, or a linear system, which has higher per-iteration complexity than four methods we used in this example.

4.3 The ℓ_1 -regularized least absolute deviation problem (LAD)

Our third example is the ℓ_1 -regularized least absolute deviation regression problem, also known as LAD in the literature. It is known that when the noise has a heavy tailed distribution such as Laplace distribution, LAD is more robust to outliers [62]. The optimization model of this problem is

$$\min_{x \in \mathbb{R}^p} \left\{ F(x) := \|Ax - b\|_1 + \lambda \|x\|_1 \right\}, \tag{30}$$

where $A \in \mathbb{R}^{n \times p}$ is generated according to a normal distribution and the noise $\mathbf{n} \in \mathbb{R}^n$ is generated by Laplace(0, 1) distribution. We generate an observed vector $b := Ax^\natural + \sigma \mathbf{n}$, where $\sigma := 0.1$ and x^\natural is a s -sparse vector of ground-truth. We choose $\lambda := 1/n$ for the regularization parameter, which gives us a good recovery of x^\natural .

This problem fits our template (1) by setting $f(\cdot) := \lambda \|\cdot\|_1$ and $g(\cdot) := \|\cdot - b\|_1$. We set β_0 in ASGARD, ASGARD-restart, and ASGARD-DL as $\beta_0 := 100\|A\|$. We

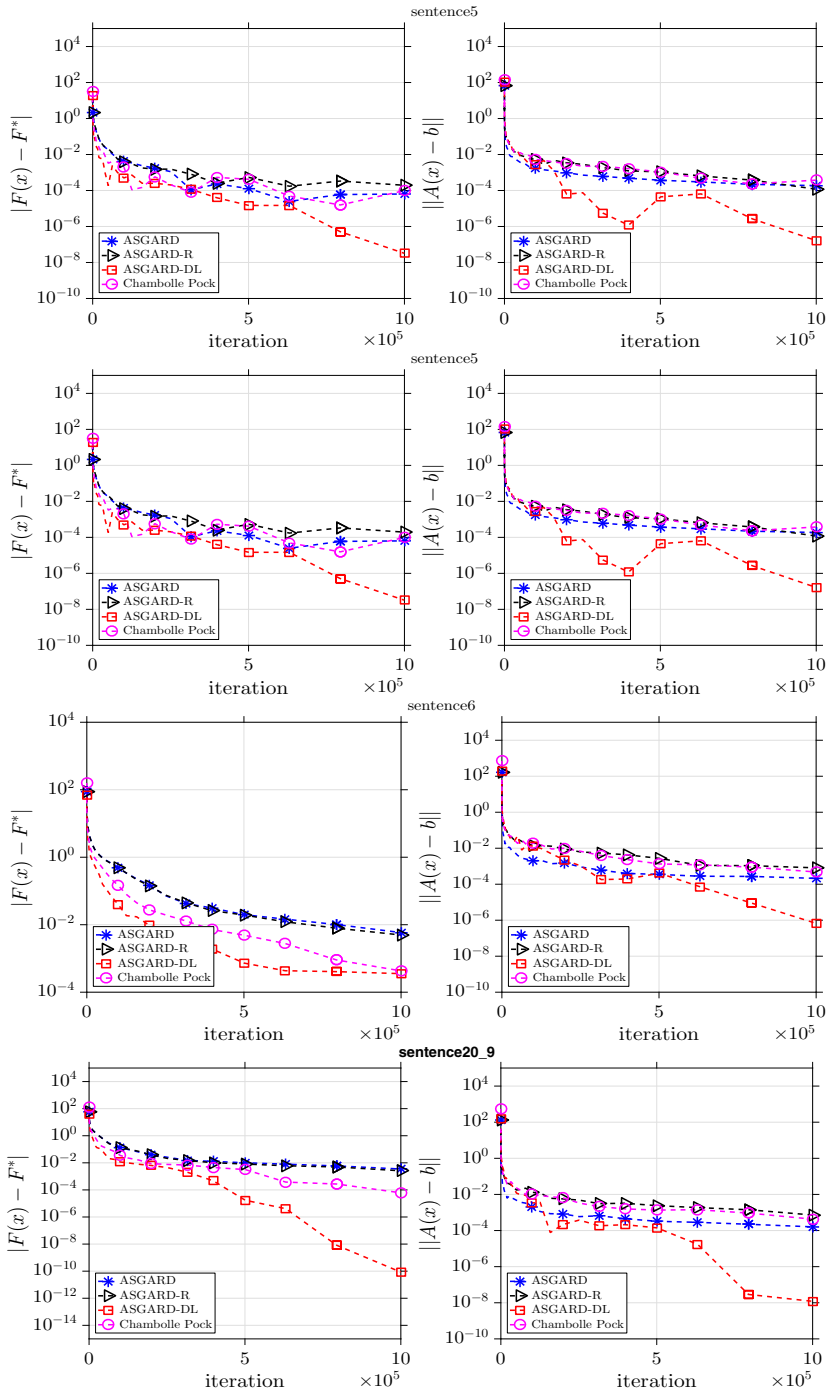


Fig. 3 The performance of 4 algorithms for solving the basis pursuit problem (29) using 4 datasets of text document

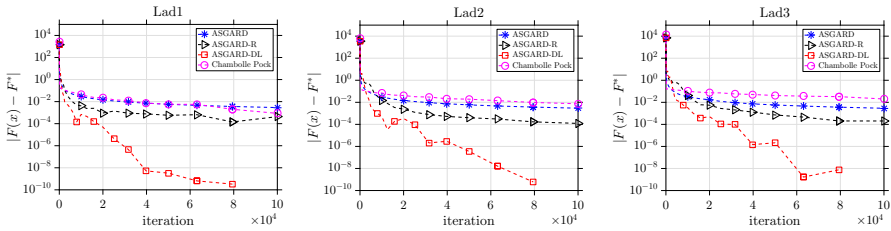


Fig. 4 The performance of 4 algorithms for solving the LAD problem (30) using three problem instances of different sizes

generated three problem instances of the size $n := 340r$, $p := 1000r$, $s := 100r$, where s is the sparsity level, and $r = 1, 2, 3$ for the first, second, and third instances, respectively. We present the results of this example in Fig. 4.

As we can see from Fig. 4 that, with the same per-iteration complexity, our method significantly outperforms the other algorithms after the accuracy of 10^{-4} . It beats other algorithms after a couple of hundred iterations and continues to decrease the objective values. Although this problem is fully nonsmooth, heuristic restart such as in ASGARD still improves the performance of the non-restart one, but does not significant outperform.

4.4 Support vector machine

Our next example is the following primal support vector machine (SVM) problem in binary classification. Instead of classical models, we consider the following ℓ_1 -regularized nonsmooth hinge loss as proposed in [69]:

$$\min_{x \in \mathbb{R}^p} \left\{ F(x) := \frac{1}{n} \sum_{i=1}^n \max \{0, 1 - b_i \langle a_i, x \rangle\} + \lambda \|x\|_1 \right\}, \tag{31}$$

where $a_i \in \mathbb{R}^p$ are the feature vectors and $b_i \in \{-1, +1\}$ are the labels for $i = 1, \dots, n$. We can cast (31) into (1) by setting $f(\cdot) := \lambda \|\cdot\|_1$ and

$$g(Ax) := \frac{1}{n} \sum_{i=1}^n \max \{0, 1 - b_i \langle a_i, x \rangle\} = \max_{y \in [0, 1]^n} \langle y, Ax + \frac{1}{n} \mathbb{1} \rangle,$$

where $A := -\frac{1}{n} [b_1 a_1, b_2 a_2, \dots, b_n a_n]^\top$ and $\mathbb{1}$ is a vector of all ones. Clearly, the proximal operator of g is simply a projection onto $[0, 1]^n$.

We use 10 different datasets from `libsvm` [15] to test four different algorithms. The initial value β_0 in ASGARD, ASGARD-restart, and ASGARD-DL is set to $\beta_0 := 0.1 \|A\|$. But for `covtype` dataset, we used $\beta_0 := 0.01 \|A\|$. The details about the datasets are given in Table 1. We test 4 algorithms on these ten datasets. The results of the first 8 problems are given in Fig. 5, and the results of the two last problems are in Fig. 6.

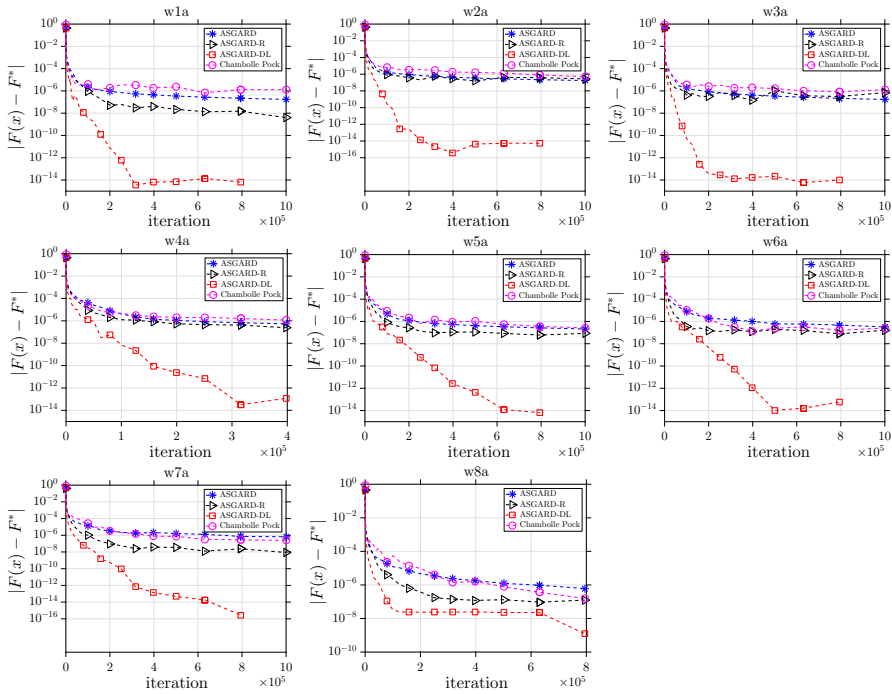


Fig. 5 The performance of 4 algorithms for solving the ℓ_1 -regularized SVM problem (31)

Table 1 The 10 datasets used for classification from libsvm [15]

Data set	Training size	Number of features
w1a	2477	300
w2a	3470	300
w3a	4912	300
w4a	7366	300
w5a	9888	300
w6a	17,188	300
w7a	24,692	300
w8a	49,749	300
rcv1	20,242	47,236
covtype	581,012	54

We again observe that Algorithm 1 significantly outperforms the other methods. Since these algorithms have the same per-iteration complexity, it is sufficient to compare them in terms of iteration numbers. Although all the algorithms have $\mathcal{O}(\frac{1}{k})$ -worst-case convergence rate, due to its double-loop, Algorithm 1 performs much better than the others, especially for high accurate solutions. This is not surprise. The double-loop allows Algorithm 1 to use large stepsize by frequently restarting τ_k and β_k , while ASGARD gradually decreases these parameters to zero, and Chambolle–Pock’s

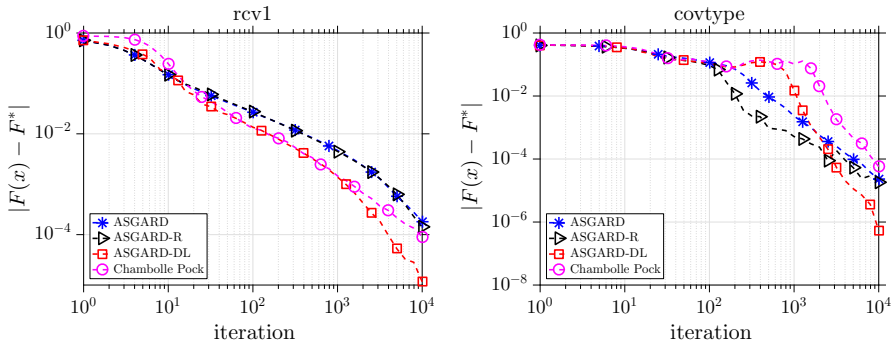


Fig. 6 The performance of 4 algorithms for solving the ℓ_1 -regularized SVM problem (31) on the two datasets {rcv1, covtype}

method fixes the step-size. Note that the $\mathcal{O}(\frac{1}{k})$ rate of Chambolle–Pock’s method is achieved via the averaged sequence, which is often much slower than the last iteration as we showed in Figs. 5 and 6.

4.5 Markowitz portfolio optimization

We consider a classical example from Markowitz portfolio optimization [11]. The setting we consider here aims at maximizing the expected return for a given risk level. Assume that we are given a vector $\rho \in \mathbb{R}^n$, where ρ is composed of expected returns from n assets. This problem can be formulated as

$$\max_{x \in \mathbb{R}^p} \left\{ \rho^\top x \mid x \in \Delta, \mathbb{E} \left[|(a_i - \rho)^\top x|^2 \right] \leq \epsilon \right\}.$$

For our setting, we use empirical sample average instead of the expectation and convert the problem to a minimization problem by negating the objective:

$$\min_{x \in \mathbb{R}^p} \left\{ -\langle \rho, x \rangle \mid x \in \Delta, \frac{1}{p} \|Ax\|_2^2 \leq \epsilon \right\}, \tag{32}$$

where $A = [(a_1 - \rho), (a_2 - \rho), \dots, (a_n - \rho)]^\top$. We map this problem to our template (26) by mapping $f(\cdot) := \delta_\Delta(\cdot)$, $g(\cdot) := \delta_{\{\|\cdot\|_2 \leq \sqrt{p\epsilon}\}}(\cdot)$, and $h(x) := -\langle \rho, x \rangle$. One key step of primal-dual algorithms is computing the projection onto an ℓ_2 -norm ball and on a simplex. Here, the complexity of simplex projection is $\mathcal{O}(p \log p)$.

As before, we apply 4 algorithms to solve (32). We use 4 datasets that are also considered in [6]. The details about the datasets are given in Table 2.

We summarized the parameters that we used for these algorithms in Table 2, where β_0 is common to ASgard, ASgard-R, and our algorithm, restart frequency (RF) is specific to ASgard-restart, ω and m_s are specific to our algorithm, and τ and σ are specific to Chambolle–Pock’s algorithm.

We have tested 4 algorithms on 4 real datasets and the results are shown in Fig. 7. As can be seen, except for the SP500 dataset, Algorithm 1 significantly outperforms

Table 2 Portfolio optimization datasets and the choice of algorithmic parameters

The size of datasets				The parameters used in 4 algorithms					
Datasets	n	p	ϵ in (32)	β_0	RF	ω	m_s	σ	τ
DJIA	507	30	0.002	$\ A\ $	10	1.1	11	$\frac{1}{\ A\ }$	$\frac{0.9999}{\ A\ }$
NYSE	5651	36	0.02	$100\ A\ $	10	1.1	11	$\frac{1}{\ A\ }$	$\frac{0.9999}{\ A\ }$
SP500	1276	25	0.02	$100\ A\ $	10	1.2	6	$\frac{1}{\ A\ }$	$\frac{0.9999}{\ A\ }$
TSE	1258	88	0.002	$100\ A\ $	10	1.1	11	$\frac{1}{\ A\ }$	$\frac{0.9999}{\ A\ }$

the other methods and shows a much faster practical performance than $\mathcal{O}(\frac{1}{k})$ rate guarantee. For SP500 dataset, ASGARD-restart algorithm shows a comparable performance to our method. However, as discussed in [59], the effect of restarting to ASGARD method is not understood theoretically. Our algorithm theoretically preserves the best-known $\mathcal{O}(\frac{1}{k})$ rate guarantee while performing as fast as, and most of the times faster than the heuristic restarting ASGARD method.

4.6 Sparse subspace clustering

In the last example, we consider the following sparse subspace clustering problem which has broad applications in machine learning, computer vision, and image processing. This problem has been studied extensively in the literature [26,27,54]. In this problem setting, we assume that there exist n points $\{x_1, x_2, \dots, x_n\} \in \mathbb{R}^p$ lying in the union of subspaces in \mathbb{R}^p . We form a matrix $X \in \mathbb{R}^{p \times n}$ by stacking $\{x_1, x_2, \dots, x_n\}$ as the columns. With this notation, each point can be represented as

$$x_j = Xc_j + e_j, \text{ s.t. } [c_j]_j = 0 \text{ and } \mathbb{1}^\top c_j = 1.$$

where $c_j \in \mathbb{R}^n$ represents the coefficients to represent point $x_j \in \mathbb{R}^p$ as an affine combination of other points, $e_j \in \mathbb{R}^p$ is the representation error, and $\mathbb{1} \in \mathbb{R}^n$ is the vector of all ones.

This formulation can be represented compactly by stacking c_j into the j th column of a matrix C as follows:

$$X = CX \text{ s.t. } \text{diag}(C) = 0, C^\top \mathbb{1} = \mathbb{1}. \tag{33}$$

The optimization problem that we tackle in this subsection is referred to as an SSC-Lasso problem in the literature, and is written by

$$\min_{C \in \mathbb{R}^{n \times n}} \left\{ \|C\|_1 + \frac{\lambda}{2} \|X - CX\|_F^2 \mid \text{diag}(C) = 0, C^\top \mathbb{1} = \mathbb{1} \right\}. \tag{34}$$

In [26,27], ADMM is used to solve (34) and recently, Pourkamali-Anaraki and Becker [54] proposed an efficient implementation of ADMM and an application

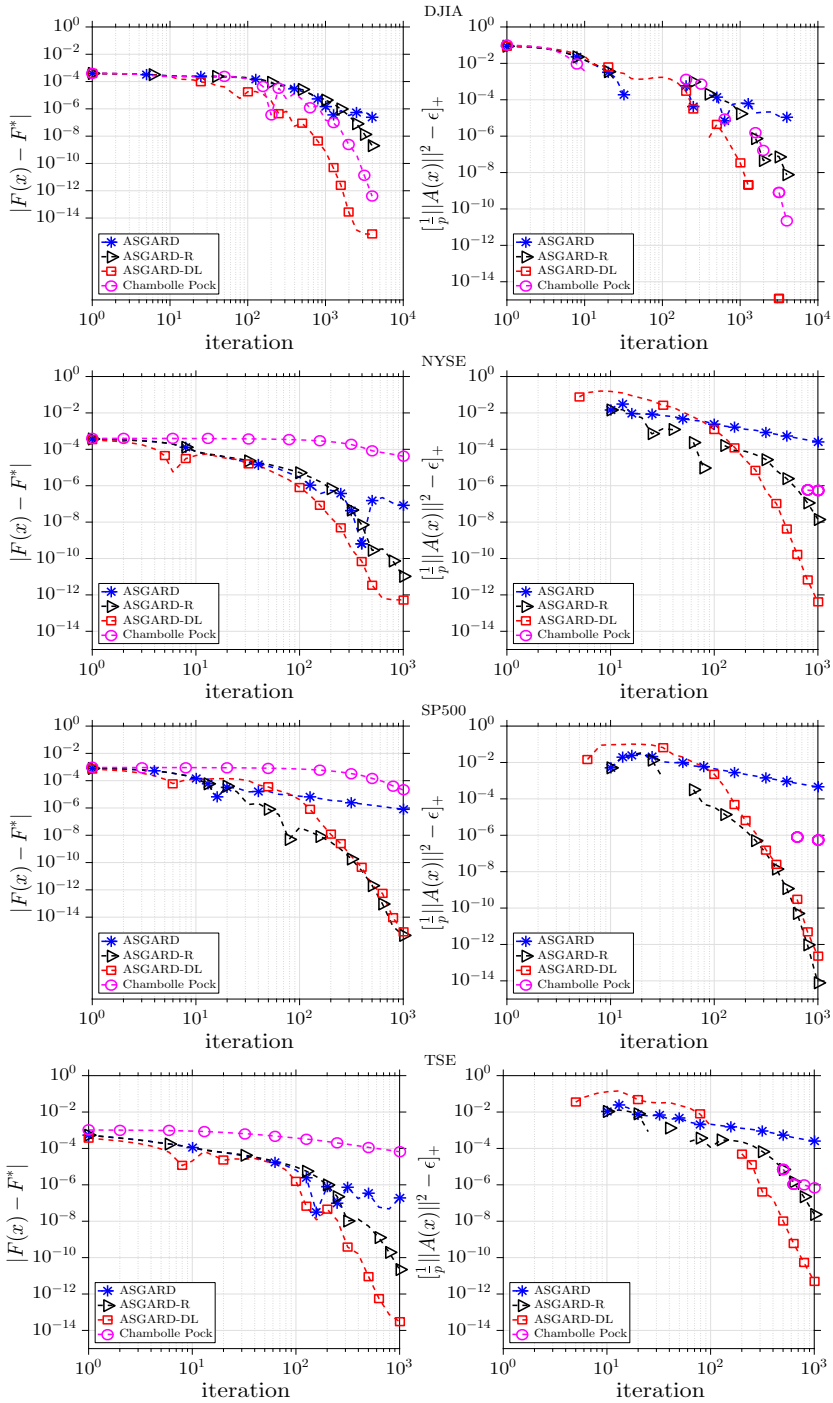


Fig. 7 The performance of 4 algorithms for solving the Markowitz portfolio optimization problem (32) on 4 real datasets

Table 3 Comparison of 3 methods on the SSC-Lasso problem (34) with $m = 2, 3, 5$ clusters and 3 independent trials of each

Problem	ADMM	TFOCS	ASGARD-DL
($n = 2$)-Objective-trial 1	236.5653	226.4371	225.7578
($n = 2$)-Clustering error-trial 1	0.0312	0.0391	0.0391
($n = 2$)-Objective-trial 2	200.3710	192.2985	191.5177
($n = 2$)-Clustering error-trial 2	0.0234	0.0469	0.0469
($n = 2$)-Objective-trial 3	197.2510	188.7655	188.1555
($n = 2$)-Clustering error-trial 3	0.0703	0.0938	0.0938
($n = 3$)-Objective-trial 1	329.9188	320.9690	319.5887
($n = 3$)-Clustering error-trial 1	0.0156	0.0156	0.0312
($n = 3$)-Objective-trial 2	341.1980	330.6395	329.5704
($n = 3$)-Clustering error-trial 2	0.0729	0.0677	0.0677
($n = 3$)-Objective-trial 3	398.8778	389.3963	388.0739
($n = 3$)-Clustering error-trial 3	0.4375	0.3594	0.3646
($n = 5$)-Objective-trial 1	549.8250	530.0340	526.1905
($n = 5$)-Clustering error-trial 1	0.1625	0.1156	0.0906
($n = 5$)-Objective-trial 2	482.8483	467.0535	461.6563
($n = 5$)-Clustering error-trial 2	0.2188	0.1125	0.1562
($n = 5$)-Objective-trial 3	1029.5459	1017.7089	1025.6752
($n = 5$)-Clustering error-trial 3	0.3156	0.3469	0.3156

of standard accelerated proximal scheme to this setting. One drawback of applying accelerated proximal schemes to (34) is the evaluation of the proximal operator of an ℓ_1 -norm over the linear constraint $C^\top \mathbb{1} = \mathbb{1}$. This requires additional computation cost of $\log(n)pn^2$. We fit (34) into our template (26) by defining $f(\cdot) := \|\cdot\|_1 + \delta_{\{\text{diag}(C)=0\}}(\cdot)$, $g(\cdot) := \delta_{\{C\mathbb{1}=\mathbb{1}\}}(\cdot)$, and $h(\cdot) = \frac{\lambda}{2}\|X - (\cdot)X\|^2$. If we apply Algorithm 1 to solve this reformulation, then no extra computation cost is incurred as in accelerated proximal gradient methods.

We use a classic benchmark Extended Yale B dataset [31] to test the sparse subspace clustering problem (34). This dataset contains face pictures of 38 individuals taken under 64 different environmental conditions. As previous works, we use downsampled images of size 48×42 pixels which correspond to $p = 2016$. We ran experiments with ADMM, TFOCS, and our method ASGARD-DL. We note that our method includes tuning parameters similar to ADMM. We use $\beta_0 := \sqrt{\|M\|}$, where $M(C) := C^\top \mathbb{1}$. We randomly selected $m = 2, 3, 5$ clusters and ran 3 trials for each case. We have used the implementation of ADMM [26,27] and TFOCS [54] provided by the authors of [26,27,54]. For fair comparison, since the per-iteration complexity of our method is smaller, we ran ADMM and TFOCS for 500 iterations and our method for 2000 iterations and saved their runtime. Then, we determined the minimum of the runtime and fetch the results for all the methods corresponding to the same computational time. We used the objective value and the clustering error as comparison measures as in [54].

Table 4 The computational time and the total number of iterations of three algorithms for solving (29) to achieve 10^{-4} and 10^{-6} accuracy levels on the KKT condition (7)

Problem	Computational time (s)			The total number of iterations		
	ASGARD-R	ASGARD-DL	CP	ASGARD-R	ASGARD-DL	CP
KKT error $\leq 10^{-4}$						
S5	44.052	27.576	8.437	144,630	93,475	29,125
S6	55.277	57.820	21.938	160,230	202,740	78,545
S19_10	52.692	51.464	37.392	171,060	172,485	126,800
S20_9	41.750	36.175	20.129	159,870	117,125	68,110
KKT error $\leq 10^{-6}$						
S5	–	122.766	240.996	–	419,005	830,005
S6	–	108.458	–	–	380,905	–
S19_10	–	104.169	–	–	353,620	–
S20_9	–	89.566	287.387	–	304,560	976,325

We can see from Table 3 that our method consistently outperforms other methods in terms of the objective values, and has similar performance in terms of the clustering errors. We present our algorithm as another candidate for solving the classical sparse subspace clustering problem with a lower per-iteration complexity than previous approaches such as ADMM and TFOCS.

4.7 Estimating moderate and high accurate KKT residuals

Our aim is to show that Algorithm 1 can potentially achieve moderate ($\varepsilon = 10^{-4}$) and high ($\varepsilon = 10^{-6}$) accuracy KKT residuals for some challenging nonsmooth instances of (1). This is not often the case in first-order primal-dual methods. We check the KKT conditions in (7) evaluated at the iterates of the algorithm. More precisely, since (7) is equivalent to $\text{prox}_{\gamma f}(x^* - \gamma A^T y^*) = 0$ and $\text{prox}_{\beta g^*}(y^* + \beta Ax^*) = 0$ for any $\gamma > 0$ and $\beta > 0$, we compute the distance $\gamma_k \|\hat{x}^{k+1} - \hat{x}^k\|$ and $\beta_s \|\tilde{y}_{k+1} - \tilde{y}^s\|$ according to Step 8 in Algorithm 1 that present an approximation of $\text{prox}_{\gamma f}(x^* - \gamma A^T y^*) = 0$ and $\text{prox}_{\beta g^*}(y^* + \beta Ax^*) = 0$, respectively. In our experiments, we used the condition $\max \{ \gamma_k \|\hat{x}^{k+1} - \hat{x}^k\|, \beta_s \|\tilde{y}_{k+1} - \tilde{y}^s\| \} \leq \varepsilon$ to terminate the algorithm. A similar condition is also used in ASGARD, ASGARD-R, and CP. Our experiment was run on a MacBook Pro. Laptop with 2.9 GHz Intel Core i7 and 16GB memory.

Firstly, we tested three algorithms: ASGARD-R, our ASGARD-DL, and CP on 4 real datasets of the basis pursuit problem (29) as an instance of the constrained setting (2). We set the maximum number of iterations in all algorithms at 10^6 . Note that we still use the same parameter settings as in Sects. 4.2 and 4.4, respectively. Since ASGARD could not achieve 10^{-4} accuracy after 10^6 iterations for any of these problems, we did not report its results. The final results are reported in Table 4, where “–” shows that the algorithm could not reach the desired accuracy after 10^6 iterations. From Table 4, we observed that CP can reach 10^{-4} accuracy faster than ASGARD-DL. ASGARD-R

can also reach 10^{-4} accuracy, but it is slower than ASGARD-DL. However, with a higher accuracy, i.e., 10^{-6} , ASGARD-DL can reach this accuracy in 4 datasets, while ASGARD-R fails in all cases and CP fails in 2 of 4 datasets. Moreover, the number of iterations in the two successful cases of CP is much higher than that of ASGARD-DL.

Next, we tested four algorithms: ASGARD, ASGARD-R, our ASGARD-DL, and CP on 10 real datasets of the support vector machine problem (31) as an instance of the composite form (1). We used the same setting as in the previous test, and the results are reported in Table 5.

In this test, ASGARD fails to reach 10^{-4} accuracy in one problem: `rcv1`, and 10^{-6} accuracy in the 7 last problems, while other three methods all achieve these accuracy levels. This is not surprise since ASGARD has $\mathcal{O}(\frac{1}{k})$ -convergence rate, and its performance is relatively tight to that bound. For 10^{-4} accuracy, ASGARD-R seems to work well and slightly outperformed ASGARD-DL and CP. However, for higher accuracy, 10^{-6} , ASGARD-R becomes slower than ASGARD-DL except for the `covtype` dataset. For `covtype`, we observed that there is an imbalance between the primal and dual optimality errors, where the primal error $\gamma_k \|\hat{x}^{k+1} - \hat{x}^k\|$ reaches 10^{-6} faster than the dual one $\beta_s \|\tilde{y}_{k+1} - \tilde{y}^s\|$. This significantly increases the number of iterations. ASGARD-DL outperforms CP in both cases: 10^{-4} and 10^{-6} accuracies except for the `covtype` dataset as we explained. Note that the KKT condition (7) used in this experiment measures the maximum error of the primal and dual optimality conditions simultaneously, while in our previous experiments, we only focused on the primal objective residual for the composite form (1), and both the primal objective residual and the primal feasibility for the constrained setting (2), but not on the dual optimality condition.

5 Further discussion and comparison with previous work

Theory and numerical methods for solving (1) and (2) are well-studied in the literature. Due to such a large proportion of solution methods, we only focus on some recent works that are most related to our method developed in this paper. We briefly survey these results to highlight the similarities and differences with our work.

In [48], Nesterov proposed combining smoothing technique and accelerated gradient methods to obtain $\mathcal{O}(\frac{1}{\varepsilon})$ -iteration complexity for attaining an ε -approximate solution to (1). However, this method requires ε to be predefined, and both the primal and dual domains are bounded. In addition, the step-size of the underlying gradient-type scheme is proportional to ε , which is often small. This leads to a poor performance in early iterations. In [47], Nesterov introduced an excessive gap technique to develop new algorithms that allow the smoothness parameters to be adaptively updated. Nevertheless, these methods still require both the primal and dual domains to be bounded, and one additional proximal operator for every two iterations.

In [13], A. Chambolle and T. Pock proposed a primal-dual algorithm to solve (1) that achieves $\mathcal{O}(\frac{1}{k})$ -convergence rate. This rate is guaranteed on a gap function and also requires both primal and dual domains to be bounded, which is unfortunately not applicable to (2). A new analysis was given in [14] to overcome this issue. Note

Table 5 The computational time and the total number of iterations of 4 algorithms for solving (31) to achieve 10^{-4} and 10^{-6} accuracy levels on the KKT condition (7)

Prob.	Computational time (s)				The total number of iterations			
	ASGD	ASGARD-R	ASGARD-DL	CP	ASGD	ASGARD-R	ASGARD-DL	CP
KKT error $\leq 10^{-4}$								
w1a	18,880	0.231	0.453	0.705	133,440	1575	4335	6445
w2a	20,524	0.305	0.505	1.360	116,095	1590	3480	9715
w3a	17,369	0.283	0.391	1.403	100,320	1410	2440	10,065
w4a	60,305	0.317	0.833	1.284	261,780	1245	4145	6990
w5a	79,265	0.409	1.029	4.085	288,290	1350	4060	17,640
w6a	208,476	0.599	1.803	4.731	491,465	1215	4430	12,525
w7a	269,159	0.736	2.250	5.652	430,840	1005	3635	8955
w8a	715,254	1.199	4.513	8.389	486,135	690	3045	5470
rcv1	-	3.777	7.660	8.417	100,000	210	1085	1355
covtype	1304,755	26.321	30.378	42.574	58,015	1170	1380	1870
KKT error $\leq 10^{-6}$								
w1a	102,919	8.403	5.496	6.942	785,915	65,280	54,915	69,325
w2a	101,371	6.464	5.244	20.529	598,780	36,945	37,290	149,995
w3a	175,270	4.888	2.532	13.587	990,150	26,070	16,380	93,940
w4a	-	11.362	6.110	28.481	-	48,345	30,900	156,100
w5a	-	14.867	10.130	31.317	-	51,795	41,005	136,775
w6a	-	15.817	10.298	79.060	-	35,715	26,030	207,400
w7a	-	27.216	13.700	187.070	-	41,205	23,035	246,765
w8a	-	84.382	40.811	453.838	-	53,730	27,930	290,575
rcv1	-	82.550	80.690	159.554	-	17,910	18,935	34,185
covtype	-	433.162	1247.115	571.864	-	19,170	56,465	27,255

that the guarantee of their methods relies on ergodic or weighted averaged sequences. In sparse and low-rank optimization and image processing, taking averaged sequence unfortunately destroys desired structures of approximate solutions. In addition, as also presented with numerical evidence, averaged sequences perform poorly in practice.

In [68], the authors proposed a homotopy algorithm called Homotopy Smoothing algorithm (HOPS) which also essentially relies on Nesterov's smoothing technique [48]. HOPS employs a similar strategy to ours in the sense of having a double loop structure. However, this method suffers from several drawbacks. First, it only applies to unconstrained problems as in (1), but not to (2) due to the unboundedness of the dual domain. Second, it requires knowing $\varepsilon_0 = P(x^0) - P(x^*)$ to be able to set the initial smoothness parameter. Third, HOPS requires tuning the number of inner iterations and the rate at which the smoothness parameter is going to be reduced. The alternative of HOPS to alleviate this issue requires a bounded primal domain which further restricts the usage of their method.

For constrained problem (2), among different methods, augmented Lagrangian schemes (ALM), alternating direction method of multipliers (ADMM), alternating minimization algorithms (AMA), and penalty methods are the most popular. Inexact augmented Lagrangian methods (iALM) [38,45,67] rely on a double loop structure similar to our method. However, the termination rules for these methods require the desired accuracy ε to be set a priori. In addition, in practice, it is not easy to check when the inner problem is solved to an ε_k -accuracy in the k -th iteration. Such an estimate is often derived from the worst-case complexity bound of the underlying solution method, and therefore, the corresponding algorithm is not efficient in practice.

While ADMM works really well and is widely used in practice, AMA is rarely used and requires additional conditions to converge. The best-known convergence rate of ADMM and its variants such as linearized ADMM and preconditioned ADMM is $\mathcal{O}(\frac{1}{k})$ under standard assumptions [34,41,42,51,66]. Moreover, this rate is given in an ergodic sense, and examples show that such a rate is optimal. See [10] for more information about the behavior of ADMM, and see [17] for the correction of a convergence issue in [10]. In practice, however, the ergodic rate is rather pessimistic, which is much slower than the last iterate sequence (see Sect. 4.1 for an example). So far, we are not aware of any work showing an $\mathcal{O}(\frac{1}{k})$ -rate of the standard ADMM or its linearized and preconditioned ADMM in the last iterate. A recent work [39] combined preconditioned/linearized ADMM and Nesterov's accelerated schemes to achieve an $\mathcal{O}(\frac{1}{k})$ -non-ergodic convergence rate.

Penalty methods use a quadratic penalty term to move the constraints to the objective and solve the subproblems by changing the penalty parameter [37,43]. Similar to iALM, these methods also do not have clear implementable termination rules for the inner loop. In addition, they do not involve dual variables. Therefore, they are often less competitive with primal-dual methods. A recent work [58] proposed a new alternating quadratic penalty algorithm to solve (2) that has the same $\mathcal{O}(\frac{1}{k})$ -non-ergodic convergence rate as in this paper. Nevertheless, this method is completely different from this paper and does not have an update on the dual center (Table 6).

Compared to our previous work [59], ASGARD, our new algorithm shares some similarities but also has several differences. First, it has inner and outer loops but the guarantee is on the overall iterations. Second, it works with any Bregman diver-

Table 6 A summary of the algorithms that require two proximal operators at each iteration

Algorithm	g is Lipschitz	$g = \delta_{\{b\}}$	Type of rate	Set ϵ	Practicality
Nesterov [48]	$P(x^k) - P^* \leq \mathcal{O}\left(\max\left(\epsilon, \frac{1}{\epsilon k^2}\right)\right)$	Not applicable	Non-ergodic	Yes	No
Chambolle-Pock [13]	$G(z^k) \leq \mathcal{O}\left(\frac{1}{k}\right)$	Convergence	Ergodic	No	No
Linearized ALM [66]	$P(z^k) - P^* \leq \mathcal{O}\left(\frac{1}{k}\right)$	$ f(z^k) - f^* \leq \mathcal{O}\left(\frac{1}{k}\right)$ $\ A_z^k - b\ \leq \mathcal{O}\left(\frac{1}{k}\right)$	Ergodic	No	No
Inexact ALM [67]	$P(z^k) - P^* \leq \mathcal{O}\left(\max(\epsilon_k, \beta_k)\right)$	$ f(z^k) - f^* \leq \mathcal{O}\left(\max(\epsilon_k + \beta_k)\right)$ $\ A_z^k - b\ \leq \mathcal{O}(\beta_k)$	Non-ergodic	Yes	No
Linearized ADMM [66]	$P(z_1^k, z_2^k) - P^* = \mathcal{O}\left(\frac{1}{k}\right)$	$ f_1(z_1^k) + f_2(z_2^k) - f_1^* - f_2^* \leq \mathcal{O}\left(\frac{1}{k}\right)$ $\ A_1 z_1^k + A_2 z_2^k - b\ \leq \mathcal{O}\left(\frac{1}{k}\right)$	Ergodic	No	No
ASGARD [59]	$P(x^k) - P^* \leq \mathcal{O}\left(\frac{1}{k}\right)$	$ f(x^k) - f^* \leq \mathcal{O}\left(\frac{1}{k}\right)$ $\ A x^k - b\ \leq \mathcal{O}\left(\frac{1}{k}\right)$	Non-ergodic	No	No
This paper (Algorithm 1)	$P(x^k) - P^* \leq \mathcal{O}\left(\frac{1}{k}\right)$	$ f(x^k) - f^* \leq \mathcal{O}\left(\frac{1}{k}\right)$ $\ A x^k - b\ \leq \mathcal{O}\left(\frac{1}{k}\right)$	Non-ergodic	No	Yes

Note that $z^k := \sum_{k=1}^K \eta_k x^k$ where K is the maximum number of iterations and η_k are the weights. For unconstrained problems, ALM/ADMM methods split the problem to obtain a constrained reformulation. The ‘‘Practicality’’ column refers to whether using the iterate in the convergence rate gives a fast practical performance or not

Table 7 A comparison with previous work [β is the smoothness parameter defined in (8)]

ADMM/iALM	Penalty/HOPS/ASGARD	This work
Constant or adaptive β	Analytically drive β to 0	Analytically drive β to 0
Update the dual center	Do not move the dual center	Update the dual center
Theory is driven by the convergence in the dual	Do not analyze the convergence of the dual	Only analyze the stability of the primal-dual sequence
Inner problems are solved inexactly	Inner problems are solved inexactly	Only ensure stability for the number of inner iterations and smoothness parameter

gence induced by a general prox-function when solving (1), while ASGARD only works with the Bregman distances induced by a strongly convex and Lipschitz gradient prox-function. This excludes some important Bregman divergences such as the Kullback–Leibler (KL) divergence. Third, our algorithm allows us to use different norms while computing proximal operators, compared to ASGARD which works with only Euclidean norms. Fourth, it automatically restarts both the primal and dual variables as well as the parameters. It also has a rigorous convergence guarantee, while the practical restarting variant of ASGARD does not have convergence guarantee.

We developed a novel analysis for our double loop structured smoothing algorithm which allowed us to derive flexible rules for parameters in both unconstrained and constrained problems, in contrast to [68]. Our analysis gives insights on the heuristic restarting strategies in [59] as well as on the number of inner iterations in the algorithm. It also gives explicit number of iterations for the inner subproblems and does not require to predefine the horizon as opposed to iALM. Table 7 summarizes the key differences between different methods we have discussed in this paper.

Acknowledgements Q. Tran-Dinh was supported in part by the National Science Foundation (NSF), Grant No. DMS-1619884, and Nafosted (no. 101.01-2017.315), Vietnam A. Alacaoglu and V. Cevher were supported in part by the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (Grant Agreement No. 725594 - time-data and Swiss National Science Foundation (SNSF) under Grant Number 200021-178865/1.

6 Appendix: The proof of technical results

This appendix provides the missing proof of the results in the main text.

6.1 The Proof of Example 1.

In this example, we have $b_{\mathcal{Y}}(y, \dot{y}) = \frac{1}{2} \|y - \dot{y}\|^2$. First, from the definition (22) of $g_{\beta}(Ax; \dot{y})$, by using the definition of $s_{\mathcal{K}}$, we write

$$g_{\beta}(Ax; \dot{y}) = \min_{u \in \mathcal{K}} \max_{y \in \mathbb{R}^n} \left\{ \langle Ax - b - u, y \rangle - \frac{\beta}{2} \|y - \dot{y}\|^2 \right\}.$$

The optimality condition of the max problem on the right hand side of the previous inequality is $Ax - b - u - \beta(y - \dot{y}) = 0$, which implies $y = \dot{y} + \frac{1}{\beta}(Ax - b - u)$. In this case, $\langle Ax - b - u, y \rangle - \frac{\beta}{2}\|y - \dot{y}\|^2 = \frac{1}{2\beta}\|Ax - b - u\|^2 + \langle \dot{y}, Ax - b - u \rangle = \frac{1}{2\beta}\|Ax - b - u + \beta\dot{y}\|^2 - \frac{\beta}{2}\|\dot{y}\|^2$. Hence, we obtain

$$\begin{aligned} g_\beta(Ax; \dot{y}) &= \min_{u \in \mathcal{K}} \left\{ \frac{1}{2\beta} \|u - (Ax - b + \beta\dot{y})\|^2 \right\} - \frac{\beta}{2} \|\dot{y}\|^2 \\ &= \frac{1}{2\beta} \text{dist}_{\mathcal{K}}(Ax - b + \beta\dot{y})^2 - \frac{\beta}{2} \|\dot{y}\|^2, \end{aligned}$$

which is (23). In addition, this implies $u = \text{proj}_{\mathcal{K}}(Ax - b + \beta\dot{y})$. Hence, we obtain $y_\beta^*(Ax; \dot{y}) = \dot{y} + \frac{1}{\beta}(Ax - b - u) = \dot{y} + \frac{1}{\beta}(Ax - b - \text{proj}_{\mathcal{K}}(Ax - b + \beta\dot{y}))$, which is exactly (24).

If \mathcal{K} is a cone, then using Moreau’s identity [2, Theorem 6.30], we can show that

$$Ax - b + \beta\dot{y} - \text{proj}_{\mathcal{K}}(Ax - b + \beta\dot{y}) = \text{proj}_{\mathcal{K}^\circ}(Ax - b + \beta\dot{y}),$$

where \mathcal{K}° is the polar set of \mathcal{K} . Since \mathcal{K} is a cone, $\mathcal{K}^\circ = -\mathcal{K}^*$, where \mathcal{K}^* is the dual cone of \mathcal{K} . Hence, we have $y_\beta^*(Ax; \dot{y}) = \text{proj}_{-\mathcal{K}^*}\left(\dot{y} + \frac{1}{\beta}(Ax - b)\right)$. □

6.2 The Proof of Lemma 1: Optimality bounds.

The proof of Lemma 1 is based on the following lemma.

Lemma 2 *Suppose that $b_{\mathcal{Y}}$ has $L_{b_{\mathcal{Y}}}$ -Lipschitz gradient such that (20) holds. Let (x^*, y^*) be a saddle point of the Lagrange function of (1) and*

$$S_\beta(\bar{x}, \dot{y}) := \max_{y \in \mathbb{R}^n} \left\{ f(\bar{x}) + \langle y, A\bar{x} \rangle - g^*(y) - \beta b_{\mathcal{Y}}(y, \dot{y}) - f(x^*) - g(Ax^*) \right\}.$$

Let $\beta_b := \beta L_{b_{\mathcal{Y}}}$ and $\bar{y}_\beta^* := y_\beta^*(A\bar{x}, \dot{y})$. Then, we have

$$\begin{cases} f(\bar{x}) + g(A\bar{x} - \beta \nabla_y b_{\mathcal{Y}}(\bar{y}_\beta^*, \dot{y})) - P^* \geq -\beta \|y^*\|_{\mathcal{Y}} \|\nabla_y b_{\mathcal{Y}}(\bar{y}_\beta^*, \dot{y})\|_{\mathcal{Y},*} \\ f(\bar{x}) + g(A\bar{x} - \beta \nabla_y b_{\mathcal{Y}}(\bar{y}_\beta^*, \dot{y})) - P^* \leq S_\beta(\bar{x}, \dot{y}) + \beta \|\dot{y}\|_{\mathcal{Y}} \|\nabla_y b_{\mathcal{Y}}(\bar{y}_\beta^*, \dot{y})\|_{\mathcal{Y},*} \\ \text{dist}_{\mathcal{Y},*}(A\bar{x}, \text{dom}(g)) \leq \beta \|\nabla_y b_{\mathcal{Y}}(\bar{y}_\beta^*, \dot{y})\|_{\mathcal{Y},*} \\ \leq \beta_b \left[\|y^* - \dot{y}\|_{\mathcal{Y}} + \left(\|y^* - \dot{y}\|_{\mathcal{Y}}^2 + \frac{2}{\beta_b} S_\beta(\bar{x}, \dot{y}) \right)^{1/2} \right]. \end{cases}$$

Proof First, from the KKT condition (7) we have $-A^\top y^* \in \partial f(x^*)$ and $y^* \in \partial g(Ax^*)$. Using these expressions and convexity of f and g , we can show that

$$\begin{aligned} f(\bar{x}) + g(A\bar{x} - \beta \nabla_y b_{\mathcal{Y}}(\bar{y}_\beta^*, \dot{y})) - f(x^*) - g(Ax^*) &\geq \langle -A^\top y^*, \bar{x} - x^* \rangle + \langle y^*, A\bar{x} - \beta \nabla_y b_{\mathcal{Y}}(\bar{y}_\beta^*, \dot{y}) - Ax^* \rangle \\ &= -\beta \langle y^*, \nabla_y b_{\mathcal{Y}}(\bar{y}_\beta^*, \dot{y}) \rangle \end{aligned}$$

$$\geq -\beta \|y^*\|_{\mathcal{Y}} \|\nabla_{\mathcal{Y}} b_{\mathcal{Y}}(\bar{y}_{\beta}^*, \dot{y})\|_{\mathcal{Y},*}, \tag{35}$$

which proves the first estimate in Lemma 2.

Next, we consider the optimality condition of

$$\bar{y}_{\beta}^* := \operatorname{argmin}_{\mathcal{Y}} \left\{ \langle A\bar{x}, y \rangle - g^*(y) - \beta b_{\mathcal{Y}}(y, \dot{y}) \right\}$$

as follows:

$$0 \in -A\bar{x} + \beta \nabla_{\mathcal{Y}} b_{\mathcal{Y}}(\bar{y}_{\beta}^*, \dot{y}) + \partial g^*(\bar{y}_{\beta}^*).$$

This implies that $\bar{y}_{\beta}^* \in \partial g(A\bar{x} - \beta \nabla_{\mathcal{Y}} b_{\mathcal{Y}}(\bar{y}_{\beta}^*, \dot{y}))$, and in particular we have $A\bar{x} - \beta \nabla_{\mathcal{Y}} b_{\mathcal{Y}}(\bar{y}_{\beta}^*, \dot{y}) \in \operatorname{dom}(g)$. Hence, we obtain

$$\operatorname{dist}_{\mathcal{Y},*}(A\bar{x}, \operatorname{dom}(g)) \leq \beta \|\nabla_{\mathcal{Y}} b_{\mathcal{Y}}(\bar{y}_{\beta}^*, \dot{y})\|_{\mathcal{Y},*}. \tag{36}$$

Now, using the equality of the Fenchel-Young inequality and the definition of $g_{\beta}(A\bar{x}, \dot{y})$, one can show that

$$\begin{aligned} & f(\bar{x}) + g(A\bar{x} - \beta \nabla_{\mathcal{Y}} b_{\mathcal{Y}}(\bar{y}_{\beta}^*, \dot{y})) - f(x^*) - g(Ax^*) \\ &= S_{\beta}(\bar{x}, \dot{y}) + g(A\bar{x} - \beta \nabla_{\mathcal{Y}} b_{\mathcal{Y}}(\bar{y}_{\beta}^*, \dot{y})) - g_{\beta}(A\bar{x}, \dot{y}) \\ &= S_{\beta}(\bar{x}, \dot{y}) + \langle \bar{y}_{\beta}^*, A\bar{x} - \beta \nabla_{\mathcal{Y}} b_{\mathcal{Y}}(\bar{y}_{\beta}^*, \dot{y}) \rangle - g^*(\bar{y}_{\beta}^*) \\ &\quad - \langle A\bar{x}, \bar{y}_{\beta}^* \rangle + g^*(\bar{y}_{\beta}^*) + \beta b_{\mathcal{Y}}(\bar{y}_{\beta}^*, \dot{y}) \\ &= S_{\beta}(\bar{x}, \dot{y}) - \beta \langle \bar{y}_{\beta}^*, \nabla_{\mathcal{Y}} b_{\mathcal{Y}}(\bar{y}_{\beta}^*, \dot{y}) \rangle + \beta b_{\mathcal{Y}}(\bar{y}_{\beta}^*, \dot{y}) \\ &\leq S_{\beta}(\bar{x}, \dot{y}) - \beta \langle \dot{y}, \nabla_{\mathcal{Y}} b_{\mathcal{Y}}(\bar{y}_{\beta}^*, \dot{y}) \rangle - \frac{\beta}{2L_{b_{\mathcal{Y}}}} \|\nabla_{\mathcal{Y}} b_{\mathcal{Y}}(\bar{y}_{\beta}^*, \dot{y})\|_{\mathcal{Y},*}^2, \end{aligned} \tag{37}$$

where we used the following bound in the last inequality of (37):

$$0 = b_{\mathcal{Y}}(\dot{y}, \dot{y}) \geq b_{\mathcal{Y}}(\bar{y}_{\beta}^*, \dot{y}) + \langle \nabla_{\mathcal{Y}} b_{\mathcal{Y}}(\bar{y}_{\beta}^*, \dot{y}), \dot{y} - \bar{y}_{\beta}^* \rangle + \frac{1}{2L_{b_{\mathcal{Y}}}} \|\nabla_{\mathcal{Y}} b_{\mathcal{Y}}(\bar{y}_{\beta}^*, \dot{y})\|_{\mathcal{Y},*}^2.$$

The second inequality of Lemma 2 follows directly from (37) using the Cauchy-Schwartz inequality.

Finally, combining (37) and (35), we obtain

$$\begin{aligned} 0 &\leq S_{\beta}(\bar{x}, \dot{y}) + \beta \langle y^* - \dot{y}, \nabla_{\mathcal{Y}} b_{\mathcal{Y}}(\bar{y}_{\beta}^*, \dot{y}) \rangle - \frac{\beta}{2L_{b_{\mathcal{Y}}}} \|\nabla_{\mathcal{Y}} b_{\mathcal{Y}}(\bar{y}_{\beta}^*, \dot{y})\|_{\mathcal{Y},*}^2 \\ &\leq S_{\beta}(\bar{x}, \dot{y}) + \beta \|y^* - \dot{y}\|_{\mathcal{Y}} \|\nabla_{\mathcal{Y}} b_{\mathcal{Y}}(\bar{y}_{\beta}^*, \dot{y})\|_{\mathcal{Y},*} - \frac{\beta}{2L_{b_{\mathcal{Y}}}} \|\nabla_{\mathcal{Y}} b_{\mathcal{Y}}(\bar{y}_{\beta}^*, \dot{y})\|_{\mathcal{Y},*}^2. \end{aligned}$$

Solving the quadratic inequation $0 \leq S_\beta(\bar{x}, \dot{y}) + \beta \|y^* - \dot{y}\|_{\mathcal{Y}} t - \frac{\beta}{2L_{b_{\mathcal{Y}}}} t^2$ in $t := \|\nabla_y b_{\mathcal{Y}}(\bar{y}_\beta^*, \dot{y})\|_{\mathcal{Y},*} \geq 0$, we deduce that

$$\|\nabla_y b_{\mathcal{Y}}(\bar{y}_\beta^*, \dot{y})\|_{\mathcal{Y},*} \leq L_{b_{\mathcal{Y}}} \left[\|y^* - \dot{y}\|_{\mathcal{Y}} + \left(\|y^* - \dot{y}\|_{\mathcal{Y}}^2 + \frac{2}{\beta L_{b_{\mathcal{Y}}}} S_\beta(\bar{x}, \dot{y}) \right)^{1/2} \right].$$

Substituting this into (36), we obtain the last estimate of Lemma 2. □

Proof of Lemma 1 The third inequality of (17) directly follows from the third inequality of Lemma 2. Note that it is void if $L_{b_{\mathcal{Y}}} = +\infty$. The first inequality of (17) is proved with the same arguments as the first inequality of Lemma 2. For the second inequality, we have two cases:

1. If $L_{b_{\mathcal{Y}}} < +\infty$, then we can use Lemma 2 to get

$$\begin{aligned} f(\bar{x}) + g(\bar{z}) - P^* &\leq f(\bar{x}) + g(A\bar{x} - \beta \nabla_y b_{\mathcal{Y}}(\bar{y}_\beta^*, \dot{y})) - P^* \\ &\quad + \widehat{M}_g \|\bar{z} - A\bar{x} + \beta \nabla_y b_{\mathcal{Y}}(\bar{y}_\beta^*, \dot{y})\|_{\mathcal{Y},*} \\ &\leq S_\beta(\bar{x}, \dot{y}) + \beta(2\widehat{M}_g + \|\dot{y}\|_{\mathcal{Y}}) \|\nabla_y b_{\mathcal{Y}}(\bar{y}_\beta^*, \dot{y})\|_{\mathcal{Y},*}, \end{aligned}$$

where we use the triangle inequality and definition of projection to get

$$\begin{aligned} \|\bar{z} - A\bar{x} + \beta \nabla_y b_{\mathcal{Y}}(\bar{y}_\beta^*, \dot{y})\|_{\mathcal{Y},*} &\leq \|\bar{z} - A\bar{x}\| + \beta \|\nabla_y b_{\mathcal{Y}}(\bar{y}_\beta^*, \dot{y})\|_{\mathcal{Y},*} \\ &= \text{dist}_{\mathcal{Y},*}(A\bar{x}, \text{dom}(g)) + \beta \|\nabla_y b_{\mathcal{Y}}(\bar{y}_\beta^*, \dot{y})\|_{\mathcal{Y},*} \\ &\stackrel{(36)}{\leq} 2\beta \|\nabla_y b_{\mathcal{Y}}(\bar{y}_\beta^*, \dot{y})\|_{\mathcal{Y},*}. \end{aligned}$$

2. If $D_{\mathcal{Y}} < +\infty$, then g is Lipschitz continuous and $\text{dom}(g) = \mathbb{R}^n$. Hence,

$$\begin{aligned} f(\bar{x}) + g(\bar{z}) - P^* &= f(\bar{x}) + g(A\bar{x}) - P^* \\ &= S_\beta(\bar{x}, \dot{y}) + g(A\bar{x}) - g_\beta(A\bar{x}, \dot{y}) \\ &\stackrel{(10)}{\leq} S_\beta(\bar{x}, \dot{y}) + \beta D_{\mathcal{Y}}. \end{aligned}$$

Combining both bounds we get the second inequality of (17). □

7 The Proof of Theorem 1: Convergence of Algorithm 1

We present the full proof of Theorem 1 in this section.

One-stage inequality: With the same argument as in [60], we can prove the following estimate at the k -th iteration at the stage s of the outer loop, i.e., $K_s \leq k < K_{s+1} := K_s + m_s$, of Algorithm 1:

$$S_{\beta_s}(\bar{x}^{k+1}; \dot{y}^s) + \frac{\tau_k^2 \|A\|^2}{\beta_s} d_{\mathcal{X}}(x^*, \hat{x}^{k+1}) \leq (1 - \tau_k) S_{\beta_s}(\bar{x}^k; \dot{y}^s) + \frac{\tau_k^2 \|A\|^2}{\beta_s} d_{\mathcal{X}}(x^*, \hat{x}^k), \tag{38}$$

where $S_{\beta}(\bar{x}; \dot{y}) := P_{\beta}(\bar{x}; \dot{y}) - P(x^*)$. Note that this estimate remains true if we use APG with Option 2. In order to use FISTA as an inner loop solver, instead of APG, we need a quadratic prox-function for the primal space. Except from this technical detail, the inequality remains the same.

Next, by strong convexity of $b_{\mathcal{Y}}(\cdot, \dot{y})$, the optimality condition of the g_{β} -subproblem, and convexity of $g^*(\cdot)$, we have

$$g_{\beta}(A\bar{x}; \dot{y}) = \max_{y \in \mathbb{R}^n} \{ \langle A\bar{x}, y \rangle - g^*(y) - \beta b_{\mathcal{Y}}(y, \dot{y}) \} \geq \langle A\bar{x}, y^* \rangle - g^*(y^*) - \beta b_{\mathcal{Y}}(y^*, \dot{y}) + \beta b_{\mathcal{Y}}(y^*, y_{\beta}^*(A\bar{x}; \dot{y})). \tag{39}$$

Now, from the KKT condition (7), we have $-A^{\top} y^* \in \partial f(x^*)$. Using this inclusion and convexity of f , we can derive

$$f(\bar{x}) \geq f(x^*) + \langle -A^{\top} y^*, \bar{x} - x^* \rangle. \tag{40}$$

Combining (39) and (40), we get

$$S_{\beta}(\bar{x}; \dot{y}) = P_{\beta}(\bar{x}; \dot{y}) - P(x^*) = f(\bar{x}) + g_{\beta}(A\bar{x}; \dot{y}) - (f(x^*) + g(Ax^*)) \geq -\beta b_{\mathcal{Y}}(y^*, \dot{y}) + \beta b_{\mathcal{Y}}(y^*, y_{\beta}^*(A\bar{x}; \dot{y})). \tag{41}$$

From (38), for $K_s \leq k \leq K_s + m_s - 1$ we obtain

$$\frac{1}{\tau_k} S_{\beta_s}(\bar{x}^{k+1}; \dot{y}^s) + \frac{\|A\|^2}{\beta_s} d_{\mathcal{X}}(x^*, \hat{x}^{k+1}) \leq \frac{1-\tau_k}{\tau_k} S_{\beta_s}(\bar{x}^k; \dot{y}^s) + \frac{\|A\|^2}{\beta_s} d_{\mathcal{X}}(x^*, \hat{x}^k). \tag{42}$$

By (41), we have $\beta b_{\mathcal{Y}}(y^*, \dot{y}) + S_{\beta}(\bar{x}; \dot{y}) \geq 0$. Let us define $D_k^s := S_{\beta_s}(\bar{x}^k; \dot{y}^s) + \beta_s b_{\mathcal{Y}}(y^*, \dot{y}^s)$. By adding $\frac{1}{\tau_k} \beta_s b_{\mathcal{Y}}(y^*, \dot{y}^s)$ to both sides of (42) and using the definition of D_k^s , we obtain

$$\frac{1}{\tau_k} D_{k+1}^s + \frac{\|A\|^2}{\beta_s} d_{\mathcal{X}}(x^*, \hat{x}^{k+1}) \leq \frac{(1-\tau_k)}{\tau_k} D_k^s + \frac{\|A\|^2}{\beta_s} d_{\mathcal{X}}(x^*, \hat{x}^k) + \frac{\beta_s}{\tau_k} b_{\mathcal{Y}}(y^*, \dot{y}^s). \tag{43}$$

Let us choose $\tau_k := \frac{2}{k - K_s + 2}$. Then, it is clear that $\tau_{K_s} = 1$. Moreover, $\frac{1-\tau_k}{\tau_k} = \frac{(k-K_s+2)(k-K_s)}{4} \leq \frac{(k-K_s+1)^2}{4} = \frac{1}{\tau_{k-1}}$. In this case, we can overestimate (43) as

$$\frac{1}{\tau_k} D_{k+1}^s + \frac{\|A\|^2}{\beta_s} d_{\mathcal{X}}(x^*, \hat{x}^{k+1}) \leq \frac{1}{\tau_{k-1}} D_k^s + \frac{\|A\|^2}{\beta_s} d_{\mathcal{X}}(x^*, \hat{x}^k) + \frac{\beta_s}{\tau_k} b_{\mathcal{Y}}(y^*, \dot{y}^s). \tag{44}$$

Taking a telescope from $k = K_s + 1$ to $k = K_{s+1} - 1 = K_s + m_s - 1$ of (44) and reusing (43) for $k = K_s$, we obtain

$$\begin{aligned}
 D_{K_{s+1}}^s + \frac{\tau_{K_{s+1}-1}^2 \|A\|^2}{\beta_s} d_{\mathcal{X}}(x^*, \hat{x}^{K_{s+1}}) &\leq \frac{\tau_{K_{s+1}-1}^2 (1 - \tau_{K_s})}{\tau_{K_s}^2} D_{K_s}^s \\
 &+ \frac{\tau_{K_{s+1}-1}^2 \|A\|^2}{\beta_s} d_{\mathcal{X}}(x^*, \hat{x}^{K_s}) + \beta_s \tau_{K_{s+1}-1}^2 b_{\mathcal{Y}}(y^*, \dot{y}^s) \sum_{j=K_s}^{K_s+m_s-1} \frac{1}{\tau_j} \\
 &\stackrel{(i)}{\leq} \frac{\tau_{K_{s+1}-1}^2 \|A\|^2}{\beta_s} d_{\mathcal{X}}(x^*, \hat{x}^{K_s}) + \beta_s \tau_{K_{s+1}-1}^2 b_{\mathcal{Y}}(y^*, \dot{y}^s) \sum_{j=K_s}^{K_s+m_s-1} \frac{1}{\tau_j},
 \end{aligned}$$

where (i) holds since $\tau_{K_s} = 1$. Since $\tau_k = \frac{2}{k - K_s + 2}$, we have $\tau_{K_{s+1}-1} = \frac{2}{m_s + 1}$ and $\sum_{j=K_s}^{K_s+m_s-1} \frac{1}{\tau_j} = \frac{m_s(m_s+3)}{4}$.

Using these relations, the last estimate leads to

$$\begin{aligned}
 D_{K_{s+1}}^s + \frac{4\|A\|^2}{(m_s+1)^2\beta_s} d_{\mathcal{X}}(x^*, \hat{x}^{K_{s+1}}) &\leq \frac{4\|A\|^2}{(m_s+1)^2\beta_s} d_{\mathcal{X}}(x^*, \hat{x}^{K_s}) \\
 &+ \frac{\beta_s m_s(m_s+3)}{(m_s+1)^2} b_{\mathcal{Y}}(y^*, \dot{y}^s). \tag{45}
 \end{aligned}$$

Since $D_{K_{s+1}}^s = S_{\beta_s}(\bar{x}^{K_{s+1}}; \dot{y}^s) + \beta_s b_{\mathcal{Y}}(y^*, \dot{y}^s) \geq \beta_s b_{\mathcal{Y}}(y^*, y_{\beta_s}^*(A\bar{x}^{K_{s+1}}; \dot{y}^s)) = \beta_s b_{\mathcal{Y}}(y^*, \dot{y}^{s+1})$ due to Step 14 of Algorithm 1, the last estimate leads to

$$\begin{aligned}
 \beta_s b_{\mathcal{Y}}(y^*, \dot{y}^{s+1}) + \frac{4\|A\|^2}{(m_s+1)^2\beta_s} d_{\mathcal{X}}(x^*, \hat{x}^{K_{s+1}}) &\leq \frac{4\|A\|^2}{(m_s+1)^2\beta_s} d_{\mathcal{X}}(x^*, \hat{x}^{K_s}) \\
 &+ \frac{\beta_s m_s(m_s+3)}{(m_s+1)^2} b_{\mathcal{Y}}(y^*, \dot{y}^s). \tag{46}
 \end{aligned}$$

The s-stage inequality. Using the update rule $m_{s+1} \leftarrow \lfloor \omega(m_s + 1) + 1 \rfloor - 1$ at Step 16 of Algorithm 1, we have

$$\omega(m_s + 1) \leq m_{s+1} + 1 \leq \omega(m_s + 1) + 1. \tag{47}$$

Hence, recursively applying this inequality, one can get

$$m_0 \omega^s \leq (m_0 + 1) \omega^s - 1 \leq m_s \leq (m_0 + \frac{\omega}{\omega-1}) \omega^s - \frac{\omega}{\omega-1} \leq \kappa_0 \omega^s, \tag{48}$$

where $\kappa_0 := m_0 + \frac{\omega}{\omega-1} > 0$.

From $\beta_{s+1} := \frac{\beta_s(m_{s+1}+1)}{\omega\sqrt{m_{s+1}(m_{s+1}+3)}}$ at Step 17 of Algorithm 1, we get $\beta_{s+1} \leq \frac{\beta_s}{\omega} \leq \frac{\beta_0}{\omega^{s+1}}$. Next, we need to lower bound β_s . Clearly, for $m_s \geq 1$, we have

$$\frac{m_{s+1}+1}{\sqrt{m_{s+1}(m_{s+1}+3)}} \geq 1 - \frac{1}{m_{s+1}} \geq 0.$$

In this case, we can estimate $\beta_{s+1} = \frac{\beta_s(m_{s+1}+1)}{\omega\sqrt{m_{s+1}(m_{s+1}+3)}} \geq \frac{\beta_s}{\omega} \left(1 - \frac{1}{m_{s+1}}\right) = \frac{\beta_s}{\omega} - \frac{\beta_s}{m_{s+1}\omega}$. Substituting (48) on m_{s+1} and β_s into this inequality, we obtain

$$\beta_{s+1} \geq \frac{\beta_s}{\omega} - \frac{c_0}{\omega^{2s+1}}, \quad \text{where } c_0 := \frac{\beta_0}{\omega m_0}.$$

This condition leads to $\omega\beta_{s+1} + \frac{c_0}{\omega^{2s}} \geq \beta_s$. By induction, we can show that $\omega^s \beta_s + c_0 \sum_{j=0}^{s-1} \frac{1}{\omega^j} \geq \beta_0$, which leads to

$$\beta_s \geq \frac{1}{\omega^s} \left(\beta_0 - \frac{c_0\omega(\omega^s-1)}{(\omega-1)\omega^s} \right) \geq \beta_0 \left(1 - \frac{1}{m_0(\omega-1)} \right) \frac{1}{\omega^s}. \tag{49}$$

Here, we use the fact that

$$\rho_0 := \beta_0 - \frac{c_0\omega(\omega^s-1)}{(\omega-1)\omega^s} \geq \beta_0 - \frac{c_0\omega}{\omega-1} = \beta_0 \left(1 - \frac{1}{m_0(\omega-1)} \right) > 0$$

since $m_0 > \frac{1}{\omega-1}$. This condition gives us a lower bound on β_s .

Now, from the update rule of β_{s+1} again, we have $\frac{\omega^2\beta_{s+1}^2 m_{s+1}(m_{s+1}+3)}{(m_{s+1}+1)^2} = \beta_s^2$ and $\frac{\omega^2}{(m_{s+1}+1)^2} \leq \frac{1}{(m_s+1)^2}$ as in (47). Plugging these estimates into (46), we obtain

$$\begin{aligned} & \frac{4\|A\|^2}{(m_{s+1}+1)^2} d\mathcal{X}(x^*, \hat{x}^{K_{s+1}}) + \frac{\beta_{s+1}^2 m_{s+1}(m_{s+1}+3)}{(m_{s+1}+1)^2} b\mathcal{Y}(y^*, \dot{y}^{s+1}) \\ & \leq \frac{1}{\omega^2} \left[\frac{4\|A\|^2}{(m_s+1)^2} d\mathcal{X}(x^*, \hat{x}^{K_s}) + \frac{\beta_s^2 m_s(m_s+3)}{(m_s+1)^2} b\mathcal{Y}(y^*, \dot{y}^s) \right]. \end{aligned} \tag{50}$$

By induction and using $\hat{x}^0 = \bar{x}^0$, we obtain

$$\begin{aligned} & \frac{4\|A\|^2}{(m_s+1)^2} d\mathcal{X}(x^*, \hat{x}^{K_s}) + \frac{\beta_s^2 m_s(m_s+3)}{(m_s+1)^2} b\mathcal{Y}(y^*, \dot{y}^s) \\ & \leq \frac{1}{\omega^{2s}} \left[\frac{4\|A\|^2}{(m_0+1)^2} d\mathcal{X}(x^*, \bar{x}^0) + \frac{\beta_0^2 m_0(m_0+3)}{(m_0+1)^2} b\mathcal{Y}(y^*, \dot{y}^0) \right]. \end{aligned} \tag{51}$$

Combining (51) and (45), we obtain

$$\begin{aligned} S_{\beta_s}(\bar{x}^{K_{s+1}}; \dot{y}^s) & \leq \frac{1}{\beta_s \omega^{2s}} \left[\frac{4\|A\|^2}{(m_0+1)^2} d\mathcal{X}(x^*, \bar{x}^0) + \frac{\beta_0^2 m_0(m_0+3)}{(m_0+1)^2} b\mathcal{Y}(y^*, \dot{y}^0) \right] \\ & \leq \frac{R_0}{\beta_s \omega^{2s}}, \end{aligned} \tag{52}$$

where $R_0 := \left[\frac{4\|A\|^2}{(m_0+1)^2} d\mathcal{X}(x^*, \bar{x}^0) + \frac{\beta_0^2 m_0(m_0+3)}{(m_0+1)^2} b\mathcal{Y}(y^*, \dot{y}^0) \right]^{1/2}$.

Using $m_s \leq \kappa_0 \omega^s$ in (48) to estimate the total number of iterations K_{s+1} of Algorithm 1 as

$$K_{s+1} = \sum_{i=0}^s m_i \leq \kappa_0 \sum_{i=0}^s \omega^i = \kappa_0 \left(\frac{\omega^{s+1}-1}{\omega-1} \right).$$

This condition leads to $\omega^s \geq \frac{(\omega-1)K_{s+1} + \kappa_0}{\omega \kappa_0}$. We use this to bound ω^s via the number of iterations K_{s+1} , and denote $\rho_0 := \beta_0 \left(1 - \frac{1}{m_0(\omega-1)} \right) > 0$.

Using (49) of β_s and m_s into (52), we obtain

$$S_{\beta_s}(\bar{x}^{K_{s+1}}; \dot{y}^s) \leq \frac{R_0^2}{\rho_0 \omega^s} \leq \frac{\omega \kappa_0 R_0^2}{\rho_0 [(\omega-1)K_{s+1} + \kappa_0]}. \tag{53}$$

Our next step is using (51) to bound $\|\dot{y}^s - y^*\|_{\mathcal{Y}}$. Clearly, $\frac{\beta_s^2 m_s (m_s + 3)}{(m_s + 1)^2} = \frac{\beta_{s-1}^2}{\omega^2} \geq \frac{\rho_0^2}{\omega^{2s}}$ by (49). Using (51), and strong convexity of $b_{\mathcal{Y}}$ with respect to the given norm, we can show that

$$\frac{1}{2} \|\dot{y}^s - y^*\|_{\mathcal{Y}}^2 \leq b_{\mathcal{Y}}(y^*, \dot{y}^s) \leq \frac{R_0^2}{\rho_0^2}. \tag{54}$$

The first estimate of (18): Let us denote \bar{z}^{K_s} the projection of $A\bar{x}^{K_s}$ onto $\text{dom}(g)$. Using Lemma 1, and we write

$$f(\bar{x}^{K_s}) + g(\bar{z}^{K_s}) - P^* \geq -\|y^*\|_{\mathcal{Y}} \text{dist}_{\mathcal{Y},*} \left(A\bar{x}^{K_s}, \text{dom}(g) \right),$$

which is the first estimate of (18).

The third estimate of (18): We define $\beta_{b,s} := \beta_s L_{b_{\mathcal{Y}}}$. Using Lemma 1, we have

$$\begin{aligned} \text{dist}_{\mathcal{Y},*} \left(A\bar{x}^{K_{s+1}}, \text{dom}(g) \right) &\leq \beta_s \|\nabla_{\mathcal{Y}} b_{\mathcal{Y}}(y_{\beta_s}^*(A\bar{x}^{K_{s+1}}, \dot{y}^s), \dot{y}^s)\|_{\mathcal{Y},*} \\ &\leq \beta_{b,s} \left[\|y^* - \dot{y}^s\|_{\mathcal{Y}} + \left(\|y^* - \dot{y}^s\|_{\mathcal{Y}}^2 + \frac{2}{\beta_{b,s}} S_{\beta_s}(\bar{x}^{K_{s+1}}, \dot{y}^s) \right)^{1/2} \right]. \end{aligned}$$

We note that, by using (49) and (52), we can bound $\frac{2S_{\beta_s}(\bar{x}^{K_{s+1}}; \dot{y}^s)}{\beta_s} \leq \frac{2R_0^2}{\rho_0^2}$. Using this upper bound and (54) we obtain

$$\begin{aligned} \text{dist}_{\mathcal{Y},*} \left(A\bar{x}^{K_s}, \text{dom}(g) \right) &\leq \beta_{b,s-1} \left(\frac{\sqrt{2}R_0}{\rho_0} + \left[\left(2 + \frac{2}{L_{b_{\mathcal{Y}}}} \right) \frac{R_0^2}{\rho_0^2} \right]^{1/2} \right) \\ &\leq \frac{(\sqrt{2} + 2)\beta_0 L_{b_{\mathcal{Y}}} \omega \kappa_0 R_0}{\rho_0 [(\omega - 1)K_s + \kappa_0]}, \end{aligned}$$

which is the third bound of (18).

The second estimate of (18): Using Lemma 1, we have

$$\begin{aligned} f(\bar{x}^{K_{s+1}}) + g(\bar{z}^{K_{s+1}}) - P^* &\leq S_{\beta_s}(\bar{x}^{K_{s+1}}, \dot{y}^s) \\ &\quad + \beta_s \min \left\{ D_{\mathcal{Y}}, (2\widehat{M}_g + \|\dot{y}^s\|_{\mathcal{Y}}) \|\nabla_{\mathcal{Y}} b_{\mathcal{Y}}(\dot{y}^{s+1}, \dot{y}^s)\|_{\mathcal{Y},*} \right\} \end{aligned}$$

$$\leq \frac{\omega\kappa_0 R_0^2}{\rho_0 [(\omega - 1)K_{s+1} + \kappa_0]} + \min \left\{ D_{\mathcal{Y}}, \frac{(\sqrt{2} + 2)L_{b_{\mathcal{Y}}}R_0}{\rho_0} \left[2\widehat{M}_g + \|y^*\|_{\mathcal{Y}} + \frac{\sqrt{2}R_0}{\rho_0} \right] \right\} \frac{\beta_0\omega\kappa_0}{(\omega - 1)K_{s+1} + \kappa_0},$$

which implies the second estimate of (18) by substituting $s + 1$ by s . \square

References

1. Arora, S., Khodak, M., Saunshi, N., Vodrahalli, K.: A compressed sensing view of unsupervised text embeddings, bag-of- n -grams, and LSTMs. In: International Conference on Learning Representations (2018)
2. Bauschke, H.H., Combettes, P.: Convex Analysis and Monotone Operators Theory in Hilbert Spaces, 2nd edn. Springer, Berlin (2017)
3. Beck, A., Teboulle, M.: A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM J. Imaging Sci.* **2**(1), 183–202 (2009)
4. Beck, A., Teboulle, M.: Smoothing and first order methods: a unified framework. *SIAM J. Optim.* **22**(2), 557–580 (2012)
5. Bertsekas, D.P.: Constrained Optimization and Lagrange Multiplier Methods. Athena Scientific, Belmont (1996)
6. Borodin, A., El-Yaniv, R., Gogan, V.: Can we learn to beat the best stock. *J. Artif. Intell. Res. (JAIR)* **21**, 579–594 (2004)
7. Borwein, J.M., Vanderwerff, J.D., et al.: Convex Functions: Constructions, Characterizations and Counterexamples, vol. 109. Cambridge University Press, Cambridge (2010)
8. Boş, R.I., Hendrich, C.: A variable smoothing algorithm for solving convex optimization problems. *TOP* **23**(1), 124–150 (2012)
9. Bot, R.I., Hendrich, C.: A double smoothing technique for solving unconstrained nondifferentiable convex optimization problems. *Comput. Optim. Appl.* **54**(2), 239–262 (2013)
10. 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)
11. Brodie, J., Daubechies, I., De Mol, C., Giannone, D., Loris, I.: Sparse and stable Markowitz portfolios. *Proc. Natl. Acad. Sci.* **106**(30), 12267–12272 (2009)
12. Candes, E., Romberg, J., Tao, T.: Robust uncertainty principles: exact signal reconstruction from highly incomplete frequency information. *IEEE Trans. Inf. Theory* **52**(2), 489–509 (2006)
13. Chambolle, A., Pock, T.: A first-order primal-dual algorithm for convex problems with applications to imaging. *J. Math. Imaging Vis.* **40**(1), 120–145 (2011)
14. Chambolle, A., Pock, T.: On the ergodic convergence rates of a first-order primal-dual algorithm. *Math. Program.* **159**(1–2), 253–287 (2016)
15. Chang, C.-C., Lin, C.-J.: LIBSVM: a library for support vector machines. *ACM Trans. Intell. Syst. Technol.* **2**, 27:1–27:27 (2011)
16. Chen, G., Teboulle, M.: Convergence analysis of a proximal-like minimization algorithm using Bregman functions. *SIAM J. Optim.* **3**(3), 538–543 (1993)
17. Chen, L., Sun, D., Toh, K.-C.: A note on the convergence of ADMM for linearly constrained convex optimization problems. *Comput. Optim. Appl.* **66**(2), 327–343 (2017)
18. Chen, S.S., Donoho, D.L., Saunders, M.A.: Atomic decomposition by basis pursuit. *SIAM Rev.* **43**(1), 129–159 (2001)
19. Combettes, P., Pesquet, J.-C.: Signal recovery by proximal forward–backward splitting. In *Fixed-Point Algorithms for Inverse Problems in Science and Engineering*, pp. 185–212. Springer (2011)
20. Condat, L.: A primal-dual splitting method for convex optimization involving Lipschitzian, proximable and linear composite terms. *J. Optim. Theory Appl.* **158**, 460–479 (2013)
21. Davis, D.: Convergence rate analysis of the forward-Douglas–Rachford splitting scheme. *SIAM J. Optim.* **25**(3), 1760–1786 (2015)

22. Davis, D., Yin, W.: Faster convergence rates of relaxed Peaceman–Rachford and ADMM under regularity assumptions. *Math. Oper. Res.* **42**(3), 577–896 (2017)
23. Devolder, O., Glineur, F., Nesterov, Y.: Double smoothing technique for large-scale linearly constrained convex optimization. *SIAM J. Optim.* **22**(2), 702–727 (2012)
24. Donoho, D.L.: Compressed sensing. *IEEE Trans. Inf. Theory* **25**(4), 1289–1306 (2006)
25. Eckstein, J.: Nonlinear proximal point algorithms using Bregman functions, with applications to convex programming. *Math. Oper. Res.* **18**(1), 202–226 (1993)
26. Elhamifar, E., Vidal, R.: Sparse subspace clustering. In: *IEEE Conference Computer Vision and Pattern Recognition (CVPR)*, pp. 2790–2797. IEEE (2009)
27. Elhamifar, E., Vidal, R.: Sparse subspace clustering: algorithm, theory, and applications. *IEEE Trans. Pattern Anal. Mach. Intell.* **35**(11), 2765–2781 (2013)
28. Fercoq, O., Qu, Z.: Restarting accelerated gradient methods with a rough strong convexity estimate. pp. 1–23 (2016). Preprint: [arXiv:1609.07358](https://arxiv.org/abs/1609.07358)
29. Gabay, D., Mercier, B.: A dual algorithm for the solution of nonlinear variational problems via finite element approximation. *Comput. Math. Appl.* **2**(1), 17–40 (1976)
30. Gao, X., Zhang, S.-Z.: First-order algorithms for convex optimization with nonseparable objective and coupled constraints. *J. Oper. Res. Soc. China* **5**(2), 131–159 (2017)
31. Georghiades, A.S., Belhumeur, P.N., Kriegman, D.J.: From few to many: illumination cone models for face recognition under variable lighting and pose. *IEEE Trans. Pattern Anal. Mach. Intell.* **23**(6), 643–660 (2001)
32. Giselsson, P., Boyd, S.: Monotonicity and restart in fast gradient methods. In: *IEEE Conference on Decision and Control, CDC*, pp. 5058–5063. Los Angeles, USA, December (2014)
33. Glowinski, R., Marrocco, A.: Sur l’approximation, par éléments finis d’ordre un, et la résolution, par pénalisation-dualité, d’une classe de problèmes de Dirichlet non linéaires. *Rev. Française Automat. Informat. Recherche Opérationnelle, RAIRO Analyse Numérique* **9**(R–2), 41–76 (1975)
34. He, B.S., Yuan, X.M.: On the $O(1/n)$ convergence rate of the Douglas–Rachford alternating direction method. *SIAM J. Numer. Anal.* **50**, 700–709 (2012)
35. Hestenes, M.R.: Multiplier and gradient methods. *J. Optim. Theory Appl.* **4**, 303–320 (1969)
36. Kiwiel, K.C.: Proximal minimization methods with generalized Bregman functions. *SIAM J. Control Optim.* **35**(4), 1142–1168 (1997)
37. Lan, G., Monteiro, R.D.C.: Iteration complexity of first-order penalty methods for convex programming. *Math. Program.* **138**(1), 115–139 (2013)
38. Lan, G., Monteiro, R.D.C.: Iteration-complexity of first-order augmented Lagrangian methods for convex programming. *Math. Program.* **155**(1–2), 511–547 (2016)
39. Li, H., Lin, Z.: Accelerated alternating direction method of multipliers: an optimal $\mathcal{O}(1/k)$ nonergodic analysis. *J. Sci. Comput.* **39**(2), 671–699 (2019)
40. Maas, A.L., Daly, R.E., Pham, P.T., Huang, D., Ng, A.Y., Potts, C.: Learning word vectors for sentiment analysis. In: *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pp. 142–150. Association for Computational Linguistics (2011)
41. Monteiro, R.D.C., Svaiter, B.F.: Iteration-complexity of block-decomposition algorithms and the alternating direction method of multipliers. *SIAM J. Optim.* **23**(1), 475–507 (2013)
42. Monteiro, R.D.C., Svaiter, B.F.: Iteration-complexity of block-decomposition algorithms and the alternating minimization augmented Lagrangian method. *SIAM J. Optim.* **23**(1), 475–507 (2013)
43. Necoara, I., Patrascu, A., Glineur, F.: Complexity of first-order inexact Lagrangian and penalty methods for conic convex programming. *Optim. Methods Softw.* **34**(2), 305–335 (2019)
44. Necoara, I., Suykens, J.A.K.: Applications of a smoothing technique to decomposition in convex optimization. *IEEE Trans. Autom. Control* **53**(11), 2674–2679 (2008)
45. Nedelcu, V., Necoara, I., Tran-Dinh, Q.: Computational complexity of inexact gradient augmented Lagrangian methods: application to constrained MPC. *SIAM J. Optim. Control* **52**(5), 3109–3134 (2014)
46. Nemirovskii, A., Yudin, D.: *Problem Complexity and Method Efficiency in Optimization*. Wiley, Hoboken (1983)
47. Nesterov, Y.: Excessive gap technique in nonsmooth convex minimization. *SIAM J. Optim.* **16**(1), 235–249 (2005)
48. Nesterov, Y.: Smooth minimization of non-smooth functions. *Math. Program.* **103**(1), 127–152 (2005)
49. Nguyen, V.Q., Fercoq, O., Cevher, V.: Smoothing technique for nonsmooth composite minimization with linear operator (2017). ArXiv preprint [arXiv:1706.05837](https://arxiv.org/abs/1706.05837)

50. O'Donoghue, B., Candes, E.: Adaptive restart for accelerated gradient schemes. *Found. Comput. Math.* **15**, 715–732 (2015)
51. Ouyang, Y., Chen, Y., Lan, G., Pasiliao, E.J.R.: An accelerated linearized alternating direction method of multiplier. *SIAM J. Imaging Sci.* **8**(1), 644–681 (2015)
52. Parikh, N., Boyd, S.: Proximal algorithms. *Found. Trends Optim.* **1**(3), 123–231 (2013)
53. Pennington, J., Socher, R., Manning, C.: Glove: Global vectors for word representation. In: *Proceedings of 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1532–1543 (2014)
54. Pourkamali-Anaraki, F., Becker, S.: Efficient Solvers for Sparse Subspace Clustering (2018). arXiv preprint [arXiv:1804.06291](https://arxiv.org/abs/1804.06291)
55. Rockafellar, R.T.: Augmented Lagrangians and applications of the proximal point algorithm in convex programming. *Math. Oper. Res.* **1**, 97–116 (1976)
56. Shefi, R., Teboulle, M.: Rate of convergence analysis of decomposition methods based on the proximal method of multipliers for convex minimization. *SIAM J. Optim.* **24**(1), 269–297 (2014)
57. Su, W., Boyd, S., Candes, E.: A differential equation for modeling Nesterov's accelerated gradient method: theory and insights. In: *Advances in Neural Information Processing Systems (NIPS)*, pp. 2510–2518 (2014)
58. Tran-Dinh, Q.: Proximal alternating penalty algorithms for constrained convex optimization. *Comput. Optim. Appl.* **72**(1), 1–43 (2019)
59. Tran-Dinh, Q., Fercoq, O., Cevher, V.: A smooth primal-dual optimization framework for nonsmooth composite convex minimization. *SIAM J. Optim.* **28**(1), 96–134 (2018)
60. Tseng, P.: On accelerated proximal gradient methods for convex–concave optimization. Submitted to *SIAM J. Optim.*, pp 1–20 (2008)
61. Vu, C.B.: A splitting algorithm for dual monotone inclusions involving co-coercive operators. *Adv. Comput. Math.* **38**(3), 667–681 (2013)
62. Wang, H., Li, G., Jiang, G.: Robust regression shrinkage and consistent variable selection through the LAD-Lasso. *J. Bus. Econ. Stat.* **25**(3), 347–355 (2007)
63. White, L., Togneri, R., Liu, W., Bennamoun, M.: Generating bags of words from the sums of their word embeddings. In: *17th International Conference on Intelligent Text Processing and Computational Linguistics (CICLing)* (2016)
64. Woodworth, B.E., Srebro, N.: Tight complexity bounds for optimizing composite objectives. In: *Advances in Neural Information Processing Systems (NIPS)*, pp. 3639–3647 (2016)
65. Wright, S.J.: *Optimization Algorithms for Data Analysis*. IAS/Park City Mathematics Series, pp. 1–49 (2017)
66. Xu, Y.: Accelerated first-order primal-dual proximal methods for linearly constrained composite convex programming. *SIAM J. Optim.* **27**(3), 1459–1484 (2017)
67. Xu, Y.: Iteration complexity of inexact augmented Lagrangian methods for constrained convex programming. *Math. Program.* (2019). <https://doi.org/10.1007/s10107-019-01425-9>
68. Xu, Yi, Yan, Yan, Lin, Qihang, Yang, Tianbao: Homotopy smoothing for non-smooth problems with lower complexity than $O(1/\epsilon)$. In: *Advances in Neural Information Processing Systems*, pp. 1208–1216 (2016)
69. Zhu, J., Rosset, S., Tibshirani, R., Hastie, T.J.: 1-norm support vector machines. In: *Advances in Neural Information Processing Systems (NIPS)*, pp. 49–56 (2004)

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Affiliations

Quoc Tran-Dinh¹ · Ahmet Alacaoglu² · Olivier Fercoq³ · Volkan Cevher²

- ✉ Quoc Tran-Dinh
quoctd@email.unc.edu
- Ahmet Alacaoglu
ahmet.alacaoglu@epfl.ch
- Olivier Fercoq
olivier.fercoq@telecom-paristech.fr
- Volkan Cevher
volkan.cevher@epfl.ch

¹ Department of Statistics and Operations Research, The University of North Carolina at Chapel Hill (UNC), 333 Hanes Hall, CB#3260, Chapel Hill, NC 27599-3260, USA

² Laboratory for Information and Inference Systems (LIONS), École Polytechnique Fédérale de Lausanne (EPFL), 1015 Lausanne, Switzerland

³ LTCI, Télécom ParisTech, Université Paris-Saclay, 75634 Paris, France