Accelerated Primal-Dual Proximal Block Coordinate Updating Methods for Constrained Convex Optimization*

Yangyang Xu[†] Shuzhong Zhang[‡]

Abstract

Block Coordinate Update (BCU) methods enjoy low per-update computational complexity because every time only one or a few block variables would need to be updated among possibly a large number of blocks. They are also easily parallelized and thus have been particularly popular for solving problems involving large-scale dataset and/or variables. In this paper, we propose a primal-dual BCU method for solving linearly constrained convex program in multiblock variables. The method is an accelerated version of a primal-dual algorithm proposed by the authors, which applies randomization in selecting block variables to update and establishes an O(1/t) convergence rate under convexity assumption. We show that the rate can be accelerated to $O(1/t^2)$ if the objective is strongly convex. In addition, if one block variable is independent of the others in the objective, we then show that the algorithm can be modified to achieve a linear rate of convergence. The numerical experiments show that the accelerated method performs stably with a single set of parameters while the original method needs to tune the parameters for different datasets in order to achieve a comparable level of performance.

Keywords: primal-dual method, block coordinate update, alternating direction method of multipliers (ADMM), accelerated first-order method.

Mathematics Subject Classification: 90C25, 95C06, 68W20.

1 Introduction

Motivated by the need to solve large-scale optimization problems and increasing capabilities in parallel computing, block coordinate update (BCU) methods have become particularly popular in recent years due to their low per-update computational complexity, low memory requirements, and their potentials in a distributive computing environment. In the context of optimization, BCU first appeared in the form of block coordinate descent (BCD) type of algorithms which can be applied to solve unconstrained smooth problems or those with separable nonsmooth terms in the objective

^{*}This work is partly supported by NSF grant DMS-1719549 and CMMI-1462408.

[†]xuy21@rpi.edu. Department of Mathematical Sciences, Rensselaer Polytechnic Institute

[‡]zhangs@umn.edu. Department of Industrial & Systems Engineering, University of Minnesota

(possibly with separable constraints). More recently, it has been developed for solving problems with nonseparable nonsmooth terms and/or constraint in a primal-dual framework.

In this paper, we consider the following linearly constrained multi-block structured optimization model:

$$\min_{x} f(x) + \sum_{i=1}^{M} g_i(x_i), \text{ s.t. } \sum_{i=1}^{M} A_i x_i = b,$$
(1)

where x is partitioned into disjoint blocks $(x_1, x_2, ..., x_M)$, f is a smooth convex function with Lipschitz continuous gradient, and each g_i is proper closed convex and possibly non-differentiable. Note that g_i can include an indicator function of a convex set \mathcal{X}_i , and thus (1) can implicitly include certain separable block constraints in addition to the nonseparable linear constraint.

Many applications arising in statistical and machine learning, image processing, and finance can be formulated in the form of (1) including the basis pursuit [7], constrained regression [23], support vector machine in its dual form [10], portfolio optimization [28], just to name a few.

Towards finding a solution for (1), we will first present an accelerated proximal Jacobian alternating direction method of multipliers (Algorithm 1), and then we generalize it to an accelerated randomized primal-dual block coordinate update method (Algorithm 2). Assuming strong convexity on the objective function, we will establish $O(1/t^2)$ convergence rate results of the proposed algorithms by adaptively setting the parameters, where t is the total number of iterations. In addition, if further assuming smoothness and the full-rankness we then obtain linear convergence of a modified method (Algorithm 3).

1.1 Related methods

Our algorithms are closely related to randomized coordinate descent methods, primal-dual coordinate update methods, and accelerated primal-dual methods. In this subsection, let us briefly review the three classes of methods and discuss their relations to our algorithms.

Randomized coordinate descent methods

In the absence of linear constraint, Algorithm 2 specializes to randomized coordinate descent (RCD), which was first proposed in [31] for smooth problems and later generalized in [27,38] to nonsmooth problems. It was shown that RCD converges sublinearly with rate O(1/t), which can be accelerated to $O(1/t^2)$ for convex problems and achieves a linear rate for strongly convex problems. By choosing multiple block variables at each iteration, [37] proposed to parallelize the RCD method and showed the same convergence results for parallelized RCD. This is similar to setting m > 1 in Algorithm 2, allowing parallel updates on the selected x-blocks.

Primal-dual coordinate update methods

In the presence of linear constraints, coordinate descent methods may fail to converge to a solution of the problem because fixing all but one block, the selected block variable may be uniquely determined by the linear constraint. To perform coordinate update to the linearly constrained problem (1), one effective approach is to update both primal and dual variables. Under this framework, the alternating direction method of multipliers (ADMM) is one popular choice. Originally, ADMM [14,17] was proposed for solving two-block structured problems with separable objective (by setting f=0 and M=2 in (1)), for which its convergence and also convergence rate have been wellestablished (see e.g. [2,13,22,29]). However, directly extending ADMM to the multi-block setting such as (1) may fail to converge; see [6] for a divergence example of the ADMM even for solving a linear system of equations. Lots of efforts have been spent on establishing the convergence of multi-block ADMM under stronger assumptions (see e.g. [4,6,16,25,26]) such as strong convexity or orthogonality conditions on the linear constraint. Without additional assumptions, modification is necessary for the ADMM applied to multi-block problems to be convergent; see [12, 19, 20, 39] for example. Very recently, [15] proposed a randomized primal-dual coordinate (RPDC) update method, whose asynchronous parallel version was then studied in [41]. Applied to (1), RPDC is a special case of Algorithm 2 with fixed parameters. It was shown that RPDC converges with rate O(1/t) under convexity assumption. More general than solving an optimization problem, primaldual coordinate (PDC) update methods have also appeared in solving fixed-point or monotone inclusion problems [9, 34–36]. However, for these problems, the PDC methods are only shown to converge but no convergence rate estimates are known unless additional assumptions are made such as the strong monotonicity condition.

Accelerated primal-dual methods

It is possible to accelerate the rate of convergence from O(1/t) to $O(1/t^2)$ for gradient type methods. The first acceleration result was shown by Nesterov [30] for solving smooth unconstrained problems. The technique has been generalized to accelerate gradient-type methods on possibly nonsmooth convex programs [1,32]. Primal-dual methods on solving linearly constrained problems can also be accelerated by similar techniques. Under convexity assumption, the augmented Lagrangian method (ALM) is accelerated in [21] from O(1/t) convergence rate to $O(1/t^2)$ by using a similar technique as that in [1] to the multiplier update, and [40] accelerates the linearized ALM using a technique similar to that in [32]. Assuming strong convexity on the objective, [18] accelerates the ADMM method, and the assumption is weakened in [40] to assuming the strong convexity for one component of the objective function. On solving bilinear saddle-point problems, various primal-dual methods can be accelerated if either primal or dual problem is strongly convex [3,5,11]. Without strong convexity, partial acceleration is still possible in terms of the rate depending on some other quantities; see e.g. [8,33].

1.2 Contributions of this paper

We accelerate the proximal Jacobian ADMM [12] and also generalize it to an accelerated primaldual coordinate updating method for linearly constrained multi-block structured convex program, where in the objective there is a nonseparable smooth function. With parameters fixed during all iterations, the generalized method reduces to that in [15] and enjoys O(1/t) convergence rate under mere convexity assumption. By adaptively setting the parameters at different iterations, we show that the accelerated method has $O(1/t^2)$ convergence rate if the objective is strongly convex. In addition, if there is one block variable that is independent of all others in the objective (but coupled in the linear constraint) and also the corresponding component function is smooth, we modify the algorithm by treating that independent variable in a different way and establish a linear convergence result. Numerically, we test the accelerated method on quadratic programming and compare it to the (nonaccelerated) RPDC method in [15]. The results demonstrate that the accelerated method performs efficiently and stably with the parameters automatically set in accordance of the analysis, while the RPDC method needs to tune its parameters for different data in order to have a comparable performance.

1.3 Nomenclature and basic facts

Notations. For a positive integer M, we denote [M] as $\{1, \ldots, M\}$. We let x_S denote the subvector of x with blocks indexed by S. Namely, if $S = \{i_1, \ldots, i_m\}$, then $x_S = (x_{i_1}, \ldots, x_{i_m})$. Similarly, A_S denotes the submatrix of A with columns indexed by S, and g_S denotes the sum of component functions indicated by S. We use $\nabla_i f(x)$ for the partial gradient of f with respect to x_i at f and f with respect to f with respect to f at f and f with respect to f and f with respect to f and f with respect to f and f are respe

$$\Delta_W(v^+, v^o, v) = \frac{1}{2} \left[\|v^+ - v\|_W^2 - \|v^o - v\|_W^2 + \|v^+ - v^o\|_W^2 \right]. \tag{2}$$

If W = I, we simply use $\Delta(v^+, v^o, v)$. Also, we denote

$$g(x) = \sum_{i=1}^{m} g_i(x_i), \quad F(x) = f(x) + g(x), \quad \Phi(\hat{x}, x, \lambda) = F(\hat{x}) - F(x) - \langle \lambda, A\hat{x} - b \rangle.$$
 (3)

Preparations. A point (x^*, λ^*) is called a Karush-Kuhn-Tucker (KKT) point of (1) if

$$0 \in \partial F(x^*) - A^{\top} \lambda^*, \quad Ax^* - b = 0. \tag{4}$$

For convex programs, the conditions in (4) are sufficient for x^* to be an optimal solution of (1), and they are also necessary if a certain qualification condition holds (e.g., the Slater condition: there is x in the interior of the domain of F such that Ax = b). Together with the convexity of F, (4) implies

$$\Phi(x, x^*, \lambda^*) \ge 0, \, \forall x. \tag{5}$$

We will use the following lemmas as basic facts. The first lemma is straightforward to verify from the definition of $\|\cdot\|_W$; the second one is similar to Lemma 3.3 in [15]; the third one is from Lemma 3.5 in [15].

Lemma 1.1 For any vectors u, v and symmetric PSD matrix W of appropriate sizes, it holds that

$$u^{\top}Wv = \frac{1}{2} \left[\|u\|_W^2 - \|u - v\|_W^2 + \|v\|_W^2 \right]. \tag{6}$$

Lemma 1.2 Given a function ϕ , for a given x and a random vector \hat{x} , if for any λ (that may depend on \hat{x}) it holds $\mathbb{E}\Phi(\hat{x}, x, \lambda) \leq \mathbb{E}\phi(\lambda)$, then for any $\gamma > 0$, we have

$$\mathbb{E}\big[F(\hat{x}) - F(x) + \gamma \|A\hat{x} - b\|\big] \le \sup_{\|\lambda\| \le \gamma} \phi(\lambda).$$

Proof. Let $\hat{\lambda} = -\frac{\gamma(A\hat{x}-b)}{\|A\hat{x}-b\|}$ if $A\hat{x}-b\neq 0$, and $\hat{\lambda}=0$ otherwise. Then

$$\Phi(\hat{x}, x, \hat{\lambda}) = F(\hat{x}) - F(x) + \gamma ||A\hat{x} - b||.$$

In addition, since $\|\hat{\lambda}\| \leq \gamma$, we have $\phi(\hat{\lambda}) \leq \sup_{\|\lambda\| \leq \gamma} \phi(\lambda)$ and thus $\mathbb{E}\phi(\hat{\lambda}) \leq \sup_{\|\lambda\| \leq \gamma} \phi(\lambda)$. Hence, we have the desired result from $\mathbb{E}\Phi(\hat{x}, x, \hat{\lambda}) \leq \mathbb{E}\phi(\hat{\lambda})$.

Lemma 1.3 Suppose $\mathbb{E}[F(\hat{x}) - F(x^*) + \gamma ||A\hat{x} - b||] \le \epsilon$. Then,

$$\mathbb{E}||A\hat{x} - b|| \le \frac{\epsilon}{\gamma - ||\lambda^*||}, \ and \ -\frac{\epsilon||\lambda^*||}{\gamma - ||\lambda^*||} \le \mathbb{E}[F(\hat{x}) - F(x^*)] \le \epsilon,$$

where (x^*, λ^*) satisfies the optimality conditions in (4), and we assume $\|\lambda^*\| < \gamma$.

Outline. The rest of the paper is organized as follows. Section 2 presents the accelerated proximal Jacobian ADMM and its convergence results. In section 3, we propose an accelerated primal-dual block coordinate update method with convergence analysis. Section 4 assumes more structure on the problem (1) and modifies the algorithm in section 3 to have linear convergence. Numerical results are provided in section 5. Finally, section 6 concludes the paper.

2 Accelerated proximal Jacobian ADMM

In this section, we propose an accelerated proximal Jacobian ADMM for solving (1). At each iteration, the algorithm updates all M block variables in parallel by minimizing a linearized proximal approximation of the augmented Lagrangian function, and then it renews the multiplier. Specifically, it iteratively performs the following updates:

$$x_i^{k+1} = \arg\min_{x_i} \left\langle \nabla_i f(x^k) - A_i^{\top} (\lambda^k - \beta_k r^k), x_i \right\rangle + g_i(x_i) + \frac{1}{2} \|x_i - x_i^k\|_{P_i^k}, \ i = 1, \dots, M,$$
 (7a)

$$\lambda^{k+1} = \lambda^k - \rho_k r^{k+1},\tag{7b}$$

where β_k and ρ_k are scalar parameters, P^k is an $M \times M$ block diagonal matrix with P_i^k as its *i*-th diagonal block for i = 1, ..., M, and $r^k = Ax^k - b$ denotes the residual. Note that (7a) consists of M independent subproblems, and they can be solved in parallel.

Algorithm 1 summarizes the proposed method. It reduces to the proximal Jacobian ADMM in [12] if β_k, ρ_k and P^k are fixed for all k and there is no nonseparable function f. We will show that adapting the parameters as the iteration progresses can accelerate the convergence of the algorithm.

Algorithm 1: Accelerated proximal Jacobian ADMM for (1)

```
1 Initialization: choose x^1, set \lambda^1 = 0, and let r^1 = Ax^1 - b
2 for k = 1, 2, ... do
3 Choose parameters \beta_k, \rho_k and a block diagonal matrix P^k
4 Let x^{k+1} \leftarrow (7a) and \lambda^{k+1} \leftarrow (7b) with r^{k+1} = Ax^{k+1} - b.
5 If a certain stopping criterion satisfied then
6 Return (x^{k+1}, \lambda^{k+1}).
```

2.1 Technical assumptions

Throughout the analysis in this section, we make the following assumptions.

Assumption 1 There exists (x^*, λ^*) satisfying the KKT conditions in (4).

Assumption 2 ∇f is Lipschitz continuous with modulus L_f .

Assumption 3 The function q is strongly convex with modulus $\mu > 0$.

The first two assumptions are standard, and the third one is for showing convergence rate of $O(1/t^2)$, where t is the number of iterations. Note that if f is strongly convex with modulus $\mu_f > 0$, we can let $f \leftarrow f - \frac{\mu_f}{2} \| \cdot \|^2$ and $g \leftarrow g + \frac{\mu_f}{2} \| \cdot \|^2$. This way, we have a convex function f and a strongly convex function g. Hence, Assumption 3 is without loss of generality. With only convexity, Algorithm 1 can be shown to converge at the rate O(1/t) with parameters fixed for all iterations, and the order 1/t is optimal as shown in the very recent work [24].

2.2 Convergence results

In this subsection, we show the $O(1/t^2)$ convergence rate result of Algorithm 1. First, we establish a result of running one iteration of Algorithm 1.

Lemma 2.1 (One-iteration analysis) Under Assumptions 2 and 3, let $\{(x^k, \lambda^k)\}$ be the sequence generated from Algorithm 1. Then for any k and (x, λ) such that Ax = b, it holds that

$$\Phi(x^{k+1}, x, \lambda)
\leq \frac{1}{2\rho_k} \left[\|\lambda - \lambda^k\|^2 - \|\lambda - \lambda^{k+1}\|^2 + \|\lambda^k - \lambda^{k+1}\|^2 \right] - \beta_k \|r^{k+1}\|^2
- \frac{1}{2} \left[\|x^{k+1} - x\|_{P^k - \beta_k A^\top A + \mu I}^2 - \|x^k - x\|_{P^k - \beta_k A^\top A}^2 + \|x^{k+1} - x^k\|_{P^k - \beta_k A^\top A - L_f I}^2 \right].$$
(8)

Using the above lemma, we are able to prove the following theorem.

Theorem 2.2 Under Assumptions 2 and 3, let $\{(x^k, \lambda^k)\}$ be the sequence generated by Algorithm 1. Suppose that the parameters are set to satisfy

$$0 < \rho_k \le 2\beta_k, \quad P^k \succeq \beta_k A^{\mathsf{T}} A + L_f I, \, \forall k \ge 1, \tag{9}$$

and there exists a number k_0 such that for all $k \geq 2$,

$$\frac{k + k_0 + 1}{\rho_k} \le \frac{k + k_0}{\rho_{k-1}},\tag{10}$$

$$(k+k_0+1)(P^k-\beta_k A^{\top} A) \leq (k+k_0)(P^{k-1}-\beta_{k-1} A^{\top} A + \mu I).$$
(11)

Then, for any (x, λ) satisfying Ax = b, we have

$$\sum_{k=1}^{t} (k+k_0+1)\Phi(x^{k+1},x,\lambda) + \sum_{k=1}^{t} \frac{k+k_0+1}{2} (2\beta_k - \rho_k) ||r^{k+1}||^2 + \frac{t+k_0+1}{2} ||x^{t+1} - x||_{P^t - \beta_t A^\top A + \mu I}^2 \leq \phi_1(x,\lambda), \tag{12}$$

where

$$\phi_1(x,\lambda) = \frac{k_0 + 2}{2\rho_1} \|\lambda - \lambda^1\|^2 + \frac{k_0 + 2}{2} \|x^1 - x\|_{P^1 - \beta_1 A^\top A}^2.$$
(13)

In the next theorem, we provide a set of parameters that satisfy the conditions in Theorem 2.2 and establish the $O(1/t^2)$ convergence rate result.

Theorem 2.3 (Convergence rate of order $1/t^2$) Under Assumptions 1 through 3, let $\{(x^k, \lambda^k)\}$ be the sequence generated by Algorithm 1 with parameters set to:

$$\beta_k = \rho_k = k\beta, \quad P^k = kP + L_f I, \forall k \ge 1,$$
 (14)

where P is a block diagonal matrix satisfying $0 \prec P - \beta A^{\top} A \leq \frac{\mu}{2} I$. Then,

$$\max\left\{\beta \|r^{t+1}\|^2, \|x^{t+1} - x^*\|_{P - \beta A^\top A}^2\right\} \le \frac{2}{t(t + k_0 + 1)}\phi_1(x^*, \lambda^*),\tag{15}$$

where $k_0 = \frac{2L_f}{\mu}$, and ϕ_1 is defined in (13). In addition, letting $\gamma = \max\{2\|\lambda^*\|, 1 + \|\lambda^*\|\}$ and

$$T = \frac{t(t+2k_0+3)}{2}, \quad \bar{x}^{t+1} = \frac{\sum_{k=1}^{t} (k+k_0+1)x^k}{T},$$

we have

$$|F(\bar{x}^{t+1}) - F(x^*)| \le \frac{1}{T} \max_{\|\lambda\| \le \gamma} \phi_1(x^*, \lambda),$$
 (16a)

$$||A\bar{x}^{t+1} - b|| \le \frac{1}{T \max\{1, ||\lambda^*||\}} \max_{||\lambda|| \le \gamma} \phi_1(x^*, \lambda).$$
(16b)

3 Accelerating randomized primal-dual block coordinate updates

In this section, we generalize Algorithm 1 to a randomized setting where the user may choose to update a subset of blocks at each iteration. Instead of updating all M block variables, we randomly choose a subset of them to renew at each iteration. Depending on the number of processors (nodes, or cores), we can choose a single or multiple block variables for each update.

3.1 The algorithm

Our algorithm is an accelerated version of the randomized primal-dual coordinate update method recently proposed in [15], for which we shall use RPDC as its acronym.¹ At each iteration, it performs a block proximal gradient update to a subset of randomly selected primal variables while keeping the remaining ones fixed, followed by an update to the multipliers. Specifically, at iteration k, it selects an index set $S_k \subset \{1, \ldots, M\}$ with cardinality m and performs the following updates:

$$x_{i}^{k+1} = \begin{cases} \underset{x_{i}}{\operatorname{arg\,min}} \langle \nabla_{i} f(x^{k}) - A_{i}^{\top} (\lambda^{k} - \beta_{k} r^{k}), x_{i} \rangle + g_{i}(x_{i}) + \frac{\eta_{k}}{2} ||x_{i} - x_{i}^{k}||^{2}, & \text{if } i \in S_{k}, \\ x_{i}^{k}, & \text{if } i \notin S_{k} \end{cases}$$
(17a)

$$r^{k+1} = r^k + \sum_{i \in S_k} A_i (x_i^{k+1} - x_i^k), \tag{17b}$$

$$\lambda^{k+1} = \lambda^k - \rho_k r^{k+1},\tag{17c}$$

where β_k , ρ_k and η_k are algorithm parameters, and their values will be determined later. Note that we use $\frac{\eta_k}{2}||x_i-x_i^k||^2$ in (17a) for simplicity. It can be replaced by a PSD matrix weighted norm square term as in (7a), and our convergence results still hold.

Algorithm 2 summarizes the above method. If the parameters β_k , ρ_k and η_k are fixed during all the iterations, i.e., constant parameters, the algorithm reduces to a special case of the RPDC method

¹In fact, [15] presents a more general algorithmic framework. It assumes two groups of variables, and each has multi-block structure. Our method in Algorithm 2 is an accelerated version of one special case of Algorithm 1 in [15].

in [15]. Adapting these parameters to the iterations, we will show that Algorithm 2 enjoys faster convergence rate than RPDC if the problem is strongly convex.

Algorithm 2: Accelerated randomized primal-dual block coordinate update method for (1)

- 1 Initialization: choose x^1 , set $\lambda^1 = 0$, let $r^1 = Ax^1 b$, and choose parameter m
- **2** for k = 1, 2, ... do
- **3** Select $S_k \subset \{1, 2, ..., M\}$ uniformly at random with $|S_k| = m$.
- 4 Choose parameters β_k , ρ_k and η_k .
- 5 Let $x^{k+1} \leftarrow (17a)$ and $\lambda^{k+1} \leftarrow (17c)$.
- 6 if a certain stopping criterion satisfied then
- 7 | Return (x^{k+1}, λ^{k+1}) .

3.2 Convergence results

In this subsection, we establish convergence results of Algorithm 2 under Assumptions 1 and 3, and also the following partial gradient Lipschitz continuity assumption.

Assumption 4 For any $S \subset \{1, ..., M\}$ with |S| = m, $\nabla_S f$ is Lipschitz continuous with a uniform constant L_m .

Note that if ∇f is Lipschitz continuous with constant L_f , then $L_m \leq L_f$ and $L_M = L_f$. In addition, if x^+ and x only differ on a set $S \subset [M]$ with cardinality m, then

$$f(x^{+}) \le f(x) + \langle \nabla f(x), x^{+} - x \rangle + \frac{L_{m}}{2} ||x^{+} - x||^{2}.$$
(18)

Similar to the analysis in section 2, we first establish a result of running one iteration of Algorithm 2. Throughout this section, we denote $\theta = \frac{m}{M}$.

Lemma 3.1 (One iteration analysis) Under Assumptions 3 and 4, let $\{(x^k, \lambda^k)\}$ be the sequence generated from Algorithm 2. Then for any x such that Ax = b, it holds

$$\mathbb{E}\left[\Phi(x^{k+1}, x, \lambda^{k+1}) + (\beta_k - \rho_k) \|r^{k+1}\|^2 + \frac{\mu}{2} \|x^{k+1} - x\|^2\right]$$

$$\leq (1 - \theta) \mathbb{E}\left[\Phi(x^k, x, \lambda^k) + \beta_k \|r^k\|^2 + \frac{\mu}{2} \|x^k - x\|^2\right] - \mathbb{E}\left[\Delta_{\eta_k I - \beta_k A^\top A}(x^{k+1}, x^k, x) - \frac{L_m}{2} \|x^{k+1} - x^k\|^2\right].$$
(19)

When $\mu = 0$ (i.e., (1) is convex), Algorithm 2 has O(1/t) convergence rate with fixed β_k, ρ_k, η_k . This can be shown from (19), and a similar result in slightly different form has been established in [15, Theorem 3.6]. For completeness, we provide its proof in the appendix.

Theorem 3.2 (Un-accelerated convergence) Under Assumptions 1 and 4, let $\{(x^k, \lambda^k)\}$ be the sequence generated from Algorithm 2 with $\beta_k = \beta$, $\rho_k = \rho$, $\eta_k = \eta$ for all k, satisfying

$$0 < \rho \le \theta \beta, \quad \eta \ge L_m + \beta \|A\|_2^2,$$

where $||A||_2$ denotes the spectral norm of A. Then

$$\left| \mathbb{E}[F(\bar{x}^t) - F(x^*)] \right| \le \frac{1}{1 + \theta(t - 1)} \max_{\|\lambda\| \le \gamma} \phi_2(x^*, \lambda), \tag{20a}$$

$$\mathbb{E}||A\bar{x}^t - b|| \le \frac{1}{(1 + \theta(t - 1)) \max\{1, ||\lambda^*||\}} \max_{\|\lambda\| \le \gamma} \phi_2(x^*, \lambda), \tag{20b}$$

where (x^*, λ^*) satisfies the KKT conditions in (4), $\gamma = \max\{\|2\lambda^*\|, 1 + \|\lambda^*\|\}$, and

$$\bar{x}^t = \frac{x^{t+1} + \theta \sum_{k=2}^t x^k}{1 + \theta(t-1)}, \quad \phi_2(x,\lambda) = (1-\theta) \left(F(x^1) - F(x) \right) + \frac{\eta}{2} ||x^1 - x||^2 + \frac{\theta ||\lambda||^2}{2\rho}.$$

When F is strongly convex, the above O(1/t) convergence rate can be accelerated to $O(1/t^2)$ by adaptively changing the parameters at each iteration. The following theorem is our main result. It shows an $O(1/t^2)$ convergence result under certain conditions on the parameters. Based on this theorem, we will give a set of parameters that satisfy these conditions, thus providing a specific scheme to choose the parameters.

Theorem 3.3 Under Assumptions 3 and 4, let $\{(x^k, \lambda^k)\}$ be the sequence generated from Algorithm 2 with parameters satisfying the following conditions for a certain number k_0 :

$$\theta(k+k_0+1) \geq 1, \forall k \geq 2, \tag{21a}$$

$$(\beta_{k-1} - \rho_{k-1})(k+k_0) \ge (1-\theta)(k+k_0+1)\beta_k, \forall 2 \le k \le t,$$
 (21b)

$$\frac{\theta(k+k_0+1)-1}{\rho_{k-1}} \geq \frac{\theta(k+k_0+2)-1}{\rho_k}, \forall 2 \leq k \leq t-1, \tag{21c}$$

$$\frac{\theta(t+k_0+1)-1}{\rho_{t-1}} \ge \frac{t+k_0+1}{\rho_t},\tag{21d}$$

$$\beta_k(k+k_0+1) \ge \beta_{k-1}(k+k_0), \forall k \ge 2,$$
 (21e)

$$(k+k_0+1)(\eta_k-L_m)I \succeq \beta_k(k+k_0+1)A^{\top}A, \forall k \ge 1,$$
 (21f)

$$(k+k_0)\eta_{k-1} + \mu(\theta(k+k_0+1)-1) \ge (k+k_0+1)\eta_k, \forall k \ge 2.$$
 (21g)

Then for any (x, λ) such that Ax = b, we have

$$(t+k_{0}+1)\mathbb{E}\Phi(x^{t+1},x,\lambda) + \sum_{k=2}^{t} (\theta(k+k_{0}+1)-1)\mathbb{E}\Phi(x^{k},x,\lambda)$$

$$\leq (1-\theta)(k_{0}+2)\mathbb{E}\left[\Phi(x^{1},x,\lambda^{1}) + \beta_{1}\|r^{1}\|^{2} + \frac{\mu}{2}\|x^{1}-x\|^{2}\right] + \frac{\eta_{1}(k_{0}+2)}{2}\mathbb{E}\|x^{1}-x\|^{2} + \frac{\theta(k_{0}+3)-1}{2\rho_{1}}\mathbb{E}\|\lambda^{1}-\lambda\|^{2} - \frac{t+k_{0}+1}{2}\mathbb{E}\|x^{t+1}-x\|_{(\mu+\eta_{t})I-\beta_{t}A^{T}A}^{2}.$$

$$(22)$$

Specifying the parameters that satisfy (21), we show $O(1/t^2)$ convergence rate of Algorithm 2.

Proposition 3.4 The following parameters satisfy all conditions in (21):

$$\beta_k = \frac{\mu(\theta k + 2 + \theta)}{2\rho \|A\|_2^2}, \, \forall k \ge 1, \tag{23a}$$

$$\rho_k = \begin{cases} \frac{\theta \beta_k}{(6-5\theta)}, & \text{for } 1 \le k \le t-1, \\ \frac{(t+k_0+1)\rho_{t-1}}{\theta(t+k_0+1)-1}, & \text{for } k = t \end{cases}$$
 (23b)

$$\eta_k = \rho \beta_k ||A||_2^2 + L_m, \, \forall k \ge 1,$$
(23c)

where $\rho \geq 1$ and

$$k_0 = \frac{4}{\theta} + \frac{2L_m}{\theta\mu}. (24)$$

Theorem 3.5 (Accelerated convergence) Under Assumptions 1, 3 and 4, let $\{(x^k, \lambda^k)\}$ be the sequence generated from Algorithm 2 with parameters taken as in (23). Then

$$\left| \mathbb{E}[F(\bar{x}^{t+1}) - F(x^*)] \right| \leq \frac{1}{T} \max_{\|\lambda\| \leq \gamma} \phi_3(x^*, \lambda), \quad \mathbb{E}\|A\bar{x}^{t+1} - b\| \leq \frac{1}{T \max\{1, \|\lambda^*\|\}} \max_{\|\lambda\| \leq \gamma} \phi_3(x^*, \lambda), \quad (25)$$

where $\gamma = \max\{2\|\lambda^*\|, 1 + \|\lambda^*\|\},$

$$\bar{x}^{t+1} = \frac{(t+k_0+1)x^{t+1} + \sum_{k=2}^{t} (\theta(k+k_0+1)-1)x^k}{T},$$

$$\phi_3(x,\lambda) = (1-\theta)(k_0+2) \left[F(x^1) - F(x) + \beta_1 ||x^1||^2 + \frac{\mu}{2} ||x^1 - x||^2 \right] + \frac{\eta_1(k_0+2)}{2} ||x^1 - x||^2 + \frac{\theta(k_0+3)-1}{2\rho_1} ||\lambda||^2$$

and

$$T = (t + k_0 + 1) + \sum_{k=2}^{t} (\theta(k + k_0 + 1) - 1).$$

In addition,

$$\mathbb{E}||x^{t+1} - x^*||^2 \le \frac{2\phi_3(x^*, \lambda^*)}{(t + k_0 + 1)\left(\frac{(\rho - 1)\mu}{2\rho}(\theta t + \theta + 2) + 2\mu + L_m\right)}.$$

4 Linearly convergent primal-dual method

In this section, we assume some more structure on (1) and show that a linear rate of convergence is possible. If there is no linear constraint, Algorithm 2 reduces to the RCD method proposed in [31]. It is well-known that RCD converges linearly if the objective is strongly convex. However, with the presence of linear constraints, mere strong convexity of the objective of the primal problem

only ensures the smoothness of its Lagrangian dual function, but not its strong concavity. Hence, in general, we do not expect linear convergence by only assuming strong convexity on the primal objective function. To ensure linear convergence on both the primal and dual variables, we need additional assumptions.

Throughout this section, we suppose that there is at least one block variable being absent in the nonseparable part of the objective, namely f. For convenience, we rename this block variable to be y, and the corresponding component function and constraint coefficient matrix as h and B. Specifically, we consider the following problem

$$\min_{x,y} f(x_1, \dots, x_M) + \sum_{i=1}^{M} g_i(x_i) + h(y), \text{ s.t. } \sum_{i=1}^{M} A_i x_i + By = b.$$
 (26)

One example of (26) is the problem that appears while computing a point on the central path of a convex program. Suppose we are interested in solving

$$\min_{x} f(x_1, \dots, x_M), \text{ s.t. } \sum_{i=1}^{M} A_i x_i \le b, x_i \ge 0, i = 1, \dots, M.$$
 (27)

Let $y = b - \sum_{i=1}^{M} A_i x_i$ and use the log-barrier function. We have the log-barrier approximation of (27) as follows:

$$\min_{x,y} f(x_1, \dots, x_M) - \mu \sum_{i=1}^{M} e^{\top} \log x_i - \mu e^{\top} \log y, \text{ s.t. } \sum_{i=1}^{M} A_i x_i + y = b,$$
 (28)

where e is the all-one vector. As μ decreases, the approximation becomes more accurate.

Towards a solution to (26), we modify Algorithm 2 by updating y-variable after the x-update. Since there is only a single y-block, to balance x and y updates, we do not renew y in every iteration but instead update it in probability $\theta = \frac{m}{M}$. Hence, roughly speaking, x and y variables are updated in the same frequency. The method is summarized in Algorithm 3.

4.1 Technical assumptions

In this section, we denote $z=(x,y,\lambda)$. Assume h is differentiable. Similar to (4), a point $z^*=(x^*,y^*,\lambda^*)$ is called a KKT point of (26) if

$$0 \in \partial F(x^*) - A^{\top} \lambda^*, \tag{32a}$$

$$\nabla h(y^*) - B^{\mathsf{T}} \lambda^* = 0, \tag{32b}$$

$$Ax^* + By^* - b = 0. (32c)$$

Besides Assumptions 3 and 4, we make two additional assumptions as follows.

Algorithm 3: Randomized primal-dual block coordinate update for (26)

- 1 Initialization: choose (x^1, y^1) , set $\lambda^1 = 0$, and choose parameters $\beta, \rho, \eta_x, \eta_y, m$.
- 2 Let $r^1 = Ax^1 + By^1 b$ and $\theta = \frac{m}{M}$.
- 3 for k = 1, 2, ... do
- 4 | Select index set $S_k \subset \{1, ..., M\}$ uniformly at random with $|S_k| = m$.
- **5** Keep $x_i^{k+1} = x_i^k, \forall i \notin S_k$ and update

$$x_i^{k+1} = \underset{x_i}{\arg\min} \left\langle \nabla_i f(x^k) - A_i^{\top} (\lambda^k - \beta r^k), x_i \right\rangle + g_i(x_i) + \frac{\eta_x}{2} ||x_i - x_i^k||^2, \text{ if } i \in S_k.$$
 (29)

Let
$$r^{k+\frac{1}{2}} = r^k + \sum_{i \in S_k} A_i (x_i^{k+1} - x_i^k)$$
.

In probability $1 - \theta$ keep $y^{k+1} = y^k$, and in probability θ let $y^{k+1} = \tilde{y}^{k+1}$, where

$$\tilde{y}^{k+1} = \arg\min_{y} h(y) - \left\langle B^{\top} (\lambda^k - \beta r^{k+\frac{1}{2}}), y \right\rangle + \frac{\eta_y}{2} \|y - y^k\|^2.$$
 (30)

Let
$$r^{k+1} = r^{k+\frac{1}{2}} + B(y^{k+1} - y^k)$$
.

7 Update the multiplier by

$$\lambda^{k+1} = \lambda^k - \rho r^{k+1}. (31)$$

 ${\bf if}\ a\ certain\ stopping\ criterion\ is\ satisfied\ {\bf then}$

| Return $(x^{k+1}, y^{k+1}, \lambda^{k+1})$.

Assumption 5 There exists $z^* = (x^*, y^*, \lambda^*)$ satisfying the KKT conditions in (32).

Assumption 6 The function h is strongly convex with modulus ν , and its gradient ∇h is Lipschitz continuous with constant L_h .

The strong convexity of F and h implies

$$F(x^{k+1}) - F(x^*) - \langle \tilde{\nabla} F(x^*), x^{k+1} - x^* \rangle \ge \frac{\mu}{2} ||x^{k+1} - x^*||^2,$$
 (33a)

$$\langle y^{k+1} - y^*, \nabla h(y^{k+1}) - \nabla h(y^*) \rangle \ge \nu \|y^{k+1} - y^*\|^2.$$
 (33b)

4.2 Convergence analysis

Similar to Lemma 3.1, we first establish a result of running one iteration of Algorithm 3. It can be proven by similar arguments to those showing Lemma 3.1.

Lemma 4.1 (One iteration analysis) Under Assumptions 3, 4, and 6, let $\{(x^k, y^k, \lambda^k)\}$ be the

sequence generated from Algorithm 3. Then for any k and (x, y, λ) such that Ax + By = b, it holds

$$\mathbb{E}\varphi(z^{k+1},z) + (\beta - \rho)\mathbb{E}\|r^{k+1}\|^2 + \frac{1}{\rho}\mathbb{E}\Delta(\lambda^{k+1},\lambda^k,\lambda)
+ \mathbb{E}\left[\Delta_P(x^{k+1},x^k,x) - \frac{L_m}{2}\|x^{k+1} - x^k\|^2\right] + \mathbb{E}\Delta_Q(y^{k+1},y^k,y)
\leq (1-\theta)\mathbb{E}\varphi(z^k,z) + \beta(1-\theta)\mathbb{E}\|r^k\|^2 + \frac{1-\theta}{\rho}\mathbb{E}\Delta(\lambda^k,\lambda^{k-1},\lambda)
+ \beta\mathbb{E}\langle A(x^{k+1}-x), B(y^{k+1}-y^k)\rangle + \beta(1-\theta)\mathbb{E}\langle B(y^k-y), A(x^{k+1}-x^k)\rangle.$$
(34)

where $P = \eta_x I - \beta A^{\mathsf{T}} A$, $Q = \eta_y I - \beta B^{\mathsf{T}} B$, and

$$\varphi(z^k, z) = F(x^k) - F(x) + \frac{\mu}{2} ||x^k - x||^2 + \langle y^k - y, \nabla h(y^k) \rangle - \langle \lambda, Ax^k + By^k - b \rangle.$$
 (35)

In the following, we let

$$\Psi(z^k, z^*) = F(x^k) - F(x^*) - \langle \tilde{\nabla} F(x^*), x^k - x^* \rangle + \langle y^k - y^*, \nabla h(y^k) - \nabla h(y^*) \rangle, \tag{36}$$

and

$$\psi(z^{k}, z^{*}; P, Q, \beta, \rho, c, \tau)
= (1 - \theta) \mathbb{E}\Psi(z^{k}, z^{*}) + \frac{\beta(1 - \theta)}{2} \mathbb{E}\|r^{k}\|^{2} + \frac{1}{2} \mathbb{E}\|x^{k} - x^{*}\|_{P + \mu(1 - \theta)I}^{2} + \frac{1}{2} \mathbb{E}\|y^{k} - y^{*}\|_{Q + \frac{\beta(1 - \theta)}{\tau}B^{\top}B}^{2}
+ \frac{1}{2\rho} \mathbb{E}\left[\|\lambda^{k} - \lambda^{*}\|^{2} - (1 - \theta)\|\lambda^{k - 1} - \lambda^{*}\|^{2} + \frac{1}{\theta}\|\lambda^{k} - \lambda^{k - 1}\|^{2}\right].$$
(37)

The following theorem is key to establishing linear convergence of Algorithm 3.

Theorem 4.2 Under Assumptions 3 through 6, let $\{(x^k, y^k, \lambda^k)\}$ be the sequence generated from Algorithm 3 with $\rho = \theta \beta$. Let $0 < \alpha < \theta$ and $\gamma = \max\left\{\frac{8\|A\|_2^2}{\alpha \mu}, \frac{8\|B\|_2^2}{\alpha \nu}\right\}$. Choose $\delta, \kappa \geq 0$ such that

$$2\begin{bmatrix} 1 - (1 - \theta)(1 + \delta) & (1 - \theta)(1 + \delta) \\ (1 - \theta)(1 + \delta) & \kappa - (1 - \theta)(1 + \delta) \end{bmatrix} \succeq \begin{bmatrix} \theta & 1 - \theta \\ 1 - \theta & \frac{1}{\theta} - (1 - \theta) \end{bmatrix}, \tag{38}$$

and positive numbers $\eta_x, \eta_y, c, \tau_1, \tau_2, \beta$ such that

$$P \succeq \beta (1 - \theta) \tau_2 A^{\top} A + L_m I \tag{39a}$$

$$Q \succeq 8cQ^{\top}Q + 4c\rho^{2}(1-\theta)(1+\frac{1}{\delta})B^{\top}BB^{\top}B + \beta\tau_{1}B^{\top}B.$$
 (39b)

Then it holds that

$$(1 - \alpha)\mathbb{E}\Psi(z^{k+1}, z^{*}) + \frac{1}{2}\mathbb{E}\|x^{k+1} - x^{*}\|_{P + (\frac{\alpha\mu}{2} + \mu)I - \frac{\beta}{\tau_{1}}A^{\top}A}^{\top} + \frac{1}{2}\mathbb{E}\|y^{k+1} - y^{*}\|_{Q + (\frac{3\alpha\nu}{2} - 8cL_{h}^{2})I}^{2}$$

$$+ (\frac{\beta - \rho}{2} + \frac{1}{\gamma})\mathbb{E}\|r^{k+1}\|^{2} - \left(c\rho^{2}\left(\kappa + 2(1 - \theta)(1 + \frac{1}{\delta})\right) + 2c(\beta - \rho)^{2}\right)\mathbb{E}\|B^{\top}r^{k+1}\|^{2}$$

$$+ \left(\frac{1}{2\rho} + \frac{c}{2}\sigma_{\min}(BB^{\top})\right)\mathbb{E}\left[\|\lambda^{k+1} - \lambda^{*}\|^{2} - (1 - \theta)\|\lambda^{k} - \lambda^{*}\|^{2} + \frac{1}{\theta}\|\lambda^{k+1} - \lambda^{k}\|^{2}\right]$$

$$\leq \psi(z^{k}, z^{*}; P, Q, \beta, \rho, c, \tau_{2}). \tag{40}$$

Using Theorem 4.2, a linear convergence rate of Algorithm 3 follows.

Theorem 4.3 Under Assumptions 3 through 6, let $\{(x^k, y^k, \lambda^k)\}$ be the sequence generated from Algorithm 3 with $\rho = \theta\beta$. Let $0 < \alpha < \theta$ and $\gamma = \max\left\{\frac{8\|A\|_2^2}{\alpha\mu}, \frac{8\|B\|_2^2}{\alpha\nu}\right\}$. Assume that B is full row-rank and $\max\{\|A\|_2, \|B\|_2\} \le 1$. Choose $\delta, \kappa, \eta_x, \eta_y, c, \beta, \tau_1, \tau_2$ satisfying (38) and (39), and in addition,

$$\frac{\alpha}{2}\mu + \theta\mu > \frac{\beta}{\tau_1} \tag{41a}$$

$$\frac{3\alpha\nu}{4} > 4cL_h^2 + \frac{\beta(1-\theta)}{2\tau_2} \tag{41b}$$

$$\frac{1}{\gamma} > c\rho^2 \left(\kappa + 2(1-\theta)(1+\frac{1}{\delta}) \right) + 2c(\beta - \rho)^2. \tag{41c}$$

Then

$$\psi(z^{k+1}, z^*; P, Q, \beta, \rho, c, \tau_2) \le \frac{1}{\eta} \psi(z^k, z^*; P, Q, \beta, \rho, c, \tau_2), \tag{42}$$

where

$$\eta = \min \left\{ \frac{1-\alpha}{1-\theta}, 1 + \frac{\frac{\alpha}{2}\mu + \theta\mu - \frac{\beta}{\tau_1}}{\eta_x + \mu(1-\theta)}, 1 + \frac{\frac{3\alpha\nu}{4} - 4cL_h^2 - \frac{\beta(1-\theta)}{2\tau_2}}{\frac{\eta_y}{2} + \frac{\beta(1-\theta)}{2\tau_2}}, \right. \\
\left. 1 + \frac{\frac{2}{\gamma} - 2c\rho^2 \left(\kappa + 2(1-\theta)(1+\frac{1}{\delta})\right) - 4c(\beta-\rho)^2}{\beta(1-\theta)}, 1 + c\rho\sigma_{\min}(BB^\top) \right\} > 1.$$

We finish this section by making a few remarks.

Remark 4.1 We can always rescale A, B and b without essentially altering the linear constraints. Hence, the assumption $\max\{\|A\|_2, \|B\|_2\} \le 1$ can be made without losing generality. From (42), it is easy to see that when $P \succ 0$ and $Q \succ 0$, (x^k, y^k) converges to (x^*, y^*) R-linearly in expectation. In addition, note that

$$\begin{split} &\|\lambda^{k+1} - \lambda^*\|^2 - (1-\theta)\|\lambda^k - \lambda^*\|^2 + \frac{1}{\theta}\|\lambda^{k+1} - \lambda^k\|^2 \\ &= \theta\|\lambda^{k+1} - \lambda^*\|^2 + 2(1-\theta)\langle\lambda^{k+1} - \lambda^*, \lambda^{k+1} - \lambda^k\rangle + (\frac{1}{\theta} - 1 + \theta)\|\lambda^{k+1} - \lambda^k\|^2 \\ &\geq \left(\theta - \frac{(1-\theta)^2}{\frac{1}{\theta} - 1 + \theta}\right)\|\lambda^{k+1} - \lambda^*\|^2 \\ &= \frac{\theta}{\frac{1}{\theta} - 1 + \theta}\|\lambda^{k+1} - \lambda^*\|^2. \end{split}$$

Hence, (42) also implies an R-linear convergence of λ^k to λ^* in expectation.

Remark 4.2 We give examples of parameters that satisfy the conditions required in Theorem 4.3. First consider the case of $\theta=1$, i.e., all blocks are updated at each iteration. In this case, we can choose $\delta=0, \kappa=\frac{1}{2}$ to satisfy (38) and $\eta_x=\beta\|A\|_2^2+L_f$ to satisfy (39a) and let $\alpha=\frac{1}{2}$ and $\tau_1=\frac{\beta}{\mu}$ to ensure that (41a) holds. Finally, choose $\eta_y>\left(\beta+\frac{\beta^2}{\mu}\right)\|B\|_2^2$ and c sufficiently small, and all other conditions in Theorem 4.3 are satisfied. Next consider the case of $\theta<1$. We can choose $\delta=\frac{\theta}{4(1-\theta)}$ and $\kappa=\frac{3}{\theta}+\frac{3\theta}{4}-2$ to satisfy (38), and let $\alpha=\frac{\theta}{2}, \tau_1=\frac{\beta}{\theta\mu}, \tau_2=\frac{2\beta(1-\theta)}{\nu}, \eta_x=\beta(1+(1-\theta)\tau_2)\|A\|_2^2+L_m$, and $\eta_y>\beta(1+\tau_1)\|B\|_2^2$. With such choices, all other conditions required in Theorem 4.3 hold when c is sufficiently small.

Remark 4.3 If there is only one x-block and there is no f function, then Algorithm 3 reduces to the so-called linearized ADMM. To show the linear convergence of the linearized ADMM, one scenario in [13, Theorem 3.1] assumes² the strong convexity of g and h, the smoothness of h, and the full row-rankness of B. In Theorem 4.3, we make the same assumptions, and so our result can be considered as a generalization.

5 Numerical experiments

The aim of this section is to test the practical performance of the proposed algorithms. We test Algorithm 2 on quadratic programming

$$\min_{x} F(x) = \frac{1}{2} x^{\top} Q x + c^{\top} x, \text{ s.t. } Ax = b, x \ge 0,$$
(43)

and Algorithm 3 on the log-barrier approximation of linear programming

$$\min_{x,y} c^{\top} x - e^{\top} \log x - e^{\top} \log y, \text{ s.t. } Ax + y = b, x_i \le u_i, \forall i.$$

$$\tag{44}$$

Quadratic programming. Two types of randomized implementations are considered: one with fixed parameters and the newly introduced one with adaptive parameters, which shall be called nonadaptive RPDC and adaptive RPDC respectively. Note that the former reduces to the method proposed in [15] when applied to (43). The purpose of the experiment is to test the effect of acceleration for the latter approach.

The data was generated randomly as follows. We let $Q = HDH^{\top} \in \mathbb{R}^{n \times n}$, where H is Gaussian randomly generated orthogonal matrix and D is a diagonal matrix with $d_{ii} = 1 + (i-1)\frac{L-1}{n-1}$, $i = 1, \ldots, n$. Hence, the smallest and largest singular values of Q are 1 and L respectively, and the objective of (43) is strongly convex with modulus 1. The components of c follow standard Gaussian distribution, and those of b follow uniform distribution on [0,1]. We let $A = [B,I] \in \mathbb{R}^{p \times n}$ to

²Besides the scenario that g and h are strongly convex, h is smooth, and B is of full row-rank, [13, Theorem 3.1] also shows linear convergence of the linearized ADMM under three other different scenarios.

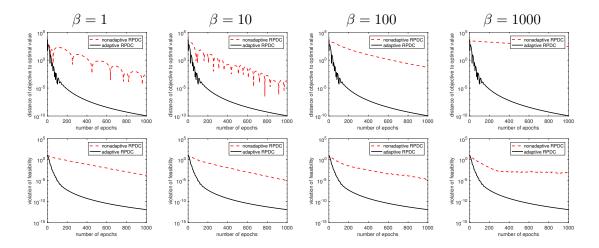


Figure 1: Results by Algorithm 2 with adaptive parameters and nonadaptive parameters for solving (43) with problem size n=2000, p=200 and condition number 10. The latter uses different penalty parameter β . Top row: difference of objective value to the optimal value $|F(x^k) - F(x^*)|$; bottom row: violation of feasibility $||Ax^k - b||$.

guarantee the existence of feasible solutions, where B was generated according to standard Gaussian distribution. In addition, we normalized A so that it has a unit spectral norm.

In the test, we fixed n=2000, p=200 and varied L among $\{10,100,1000\}$. For both nonadaptive and adaptive RPDC, we evenly partitioned x into 40 blocks, i.e., each block consists of 50 coordinates, and we set m=40, i.e., all blocks are updated at each iteration. For the adaptive RPDC, we set the values of its parameters according to (23) with $\rho=1$, and those for the nonadaptive RPDC were set based on Theorem 3.2 with $\rho=\beta, \ \eta=100+\beta, \ \forall k$ where β varied among $\{1,10,100,1000\}$. Figures 1 through 3 plot the objective values and feasibility violations by Algorithm 2 under these two different settings. From these results, we see that adaptive RPDC performed well for all three datasets with a single set of parameters while the performance of the nonadaptive one was severely affected by the penalty parameter.

Linear programming. In this test, we apply Algorithm 3 to the problem (44), where we let $f(x) = c^{\top}x$, $g(x) = -e^{\top}\log x$ and $h(y) = -e^{\top}\log y$. The purpose of this experiment is to demonstrate the linear convergence of Algorithm 3.

We generated $A \in \mathbb{R}^{200 \times 2000}$ and c according to the standard Gaussian distribution and b by the uniform distribution on $[\frac{1}{2}, \frac{3}{2}]$. The upper bound was set to $u_i = 10, \forall i$. We treated x as a single block and set the algorithm parameters to $\beta = 0.1$, $\eta_x = \beta ||A||_2^2$, and $\eta_y = \beta \left(1 + \frac{2.001\beta}{3\mu}\right)$. This setting satisfies the conditions required in Theorem 4.3 if α is sufficiently close to 1. Note that g and h do not have uniform strong convexity constants but they are both strongly convex on a bounded set. Figure 4 shows the convergence behavior of Algorithm 3. From the figure, we can clearly see that the algorithm linearly converges to an optimal solution.

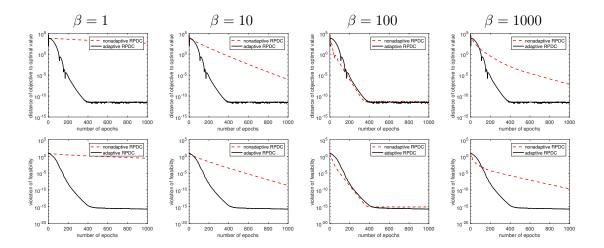


Figure 2: Results by Algorithm 2 with adaptive parameters and nonadaptive parameters for solving (43) with problem size n = 2000, p = 200 and condition number 100. The latter uses different penalty parameter β . Top row: difference of objective value to the optimal value $|F(x^k) - F(x^*)|$; bottom row: violation of feasibility $||Ax^k - b||$.

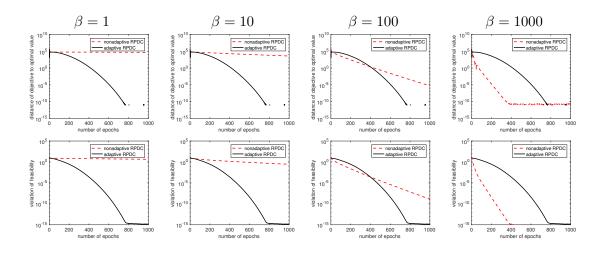


Figure 3: Results by Algorithm 2 with adaptive parameters and nonadaptive parameters for solving (43) with problem size n = 2000, p = 200 and condition number 1000. The latter uses different penalty parameter β . Top row: difference of objective value to the optimal value $|F(x^k) - F(x^*)|$; bottom row: violation of feasibility $||Ax^k - b||$.

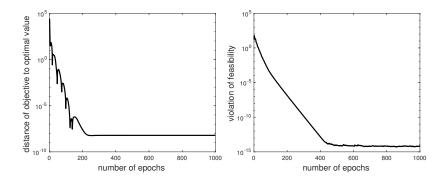


Figure 4: Results by Algorithm 3 on the problem (44) with $A \in \mathbb{R}^{200 \times 2000}$. Left: difference of objective value to the optimal value $|F(x^k) + h(y^k) - F(x^*) - h(y^*)|$; Right: violation of feasibility $||Ax^k + By^k - b||$

6 Conclusions

In this paper we propose an accelerated proximal Jacobian ADMM method and generalize it to an accelerated randomized primal-dual coordinate updating method for solving linearly constrained multi-block structured convex programs. We show that if the objective is strongly convex then the methods achieve $O(1/t^2)$ convergence rate where t is the total number of iterations. In addition, if one block variable is independent of others in the objective and its part of the objective function is smooth, we have modified the primal-dual coordinate updating method to achieve linear convergence. Numerical experiments on quadratic programming and log-barrier approximation of linear programming have shown the efficacy of the newly proposed methods.

References

- [1] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. SIAM journal on imaging sciences, 2(1):183–202, 2009. 3
- [2] D. Boley. Local linear convergence of the alternating direction method of multipliers on quadratic or linear programs. SIAM Journal on Optimization, 23(4):2183–2207, 2013. 3
- [3] K. Bredies and H. Sun. Accelerated douglas-rachford methods for the solution of convex-concave saddle-point problems. arXiv preprint arXiv:1604.06282, 2016. 3
- [4] X. Cai, D. Han, and X. Yuan. The direct extension of admm for three-block separable convex minimization models is convergent when one function is strongly convex. *Optimization Online*, 2014. 3
- [5] A. Chambolle and T. Pock. A first-order primal-dual algorithm for convex problems with applications to imaging. *Journal of Mathematical Imaging and Vision*, 40(1):120–145, 2011. 3

- [6] C. Chen, B. He, Y. Ye, and X. Yuan. The direct extension of admm for multi-block convex minimization problems is not necessarily convergent. *Mathematical Programming*, 155(1-2):57– 79, 2016. 3
- [7] S. S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. *SIAM review*, 43(1):129–159, 2001. 2
- [8] Y. Chen, G. Lan, and Y. Ouyang. Optimal primal-dual methods for a class of saddle point problems. SIAM Journal on Optimization, 24(4):1779–1814, 2014. 3
- [9] P. L. Combettes and J.-C. Pesquet. Stochastic quasi-fejér block-coordinate fixed point iterations with random sweeping. SIAM Journal on Optimization, 25(2):1221–1248, 2015. 3
- [10] C. Cortes and V. Vapnik. Support-vector networks. Machine learning, 20(3):273–297, 1995. 2
- [11] C. Dang and G. Lan. Randomized methods for saddle point computation. arXiv preprint arXiv:1409.8625, 2014. 3
- [12] W. Deng, M.-J. Lai, Z. Peng, and W. Yin. Parallel multi-block admm with o(1/k) convergence. Journal of Scientific Computing, pages 1–25, 2016. 3, 4, 6
- [13] W. Deng and W. Yin. On the global and linear convergence of the generalized alternating direction method of multipliers. *Journal of Scientific Computing*, 66(3):889–916, 2015. 3, 16
- [14] D. Gabay and B. Mercier. A dual algorithm for the solution of nonlinear variational problems via finite element approximation. *Computers & Mathematics with Applications*, 2(1):17–40, 1976. 3
- [15] X. Gao, Y. Xu, and S. Zhang. Randomized primal-dual proximal block coordinate updates. arXiv preprint arXiv:1605.05969, 2016. 3, 4, 5, 8, 9, 16
- [16] X. Gao and S.-Z. Zhang. First-order algorithms for convex optimization with nonseparable objective and coupled constraints. *Journal of the Operations Research Society of China*, pages 1–29, 2015. 3
- [17] R. Glowinski and A. Marrocco. Sur l'approximation, par eléments finis d'ordre un, et la résolution, par pénalisation-dualité d'une classe de problèmes de dirichlet non linéaires. ESAIM: Mathematical Modelling and Numerical Analysis, 9(R2):41–76, 1975. 3
- [18] T. Goldstein, B. O'Donoghue, S. Setzer, and R. Baraniuk. Fast alternating direction optimization methods. SIAM Journal on Imaging Sciences, 7(3):1588–1623, 2014. 3
- [19] B. He, L. Hou, and X. Yuan. On full Jacobian decomposition of the augmented Lagrangian method for separable convex programming. SIAM Journal on Optimization, 25(4):2274–2312, 2015. 3

- [20] B. He, M. Tao, and X. Yuan. Alternating direction method with gaussian back substitution for separable convex programming. SIAM Journal on Optimization, 22(2):313–340, 2012. 3
- [21] B. He and X. Yuan. On the acceleration of augmented lagrangian method for linearly constrained optimization. *Optimization online*, 2010. 3
- [22] B. He and X. Yuan. On the O(1/n) convergence rate of the Douglas-Rachford alternating direction method. SIAM Journal on Numerical Analysis, 50(2):700-709, 2012. 3
- [23] G. M. James, C. Paulson, and P. Rusmevichientong. Penalized and constrained regression. Technical report, 2013. 2
- [24] H. Li and Z. Lin. Optimal nonergodic o(1/k) convergence rate: When linearized adm meets nesterov's extrapolation. $arXiv\ preprint\ arXiv:1608.06366$, 2016. 6
- [25] M. Li, D. Sun, and K.-C. Toh. A convergent 3-block semi-proximal ADMM for convex minimization problems with one strongly convex block. Asia-Pacific Journal of Operational Research, 32(04):1550024, 2015. 3
- [26] T. Lin, S. Ma, and S. Zhang. On the global linear convergence of the admm with multiblock variables. SIAM Journal on Optimization, 25(3):1478–1497, 2015. 3
- [27] Z. Lu and L. Xiao. On the complexity analysis of randomized block-coordinate descent methods. *Mathematical Programming*, 152(1-2):615–642, aug 2015. 2
- [28] H. Markowitz. Portfolio selection. The journal of finance, 7(1):77–91, 1952. 2
- [29] R. D. Monteiro and B. F. Svaiter. Iteration-complexity of block-decomposition algorithms and the alternating direction method of multipliers. *SIAM Journal on Optimization*, 23(1):475–507, 2013. 3
- [30] Y. Nesterov. A method of solving a convex programming problem with convergence rate $O(1/k^2)$. Soviet Mathematics Doklady, 27(2):372–376, 1983. 3
- [31] Y. Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. SIAM Journal on Optimization, 22(2):341–362, 2012. 2, 11
- [32] Y. Nesterov. Gradient methods for minimizing composite functions. *Mathematical Programming*, 140(1):125–161, 2013. 3
- [33] Y. Ouyang, Y. Chen, G. Lan, and E. Pasiliao Jr. An accelerated linearized alternating direction method of multipliers. SIAM Journal on Imaging Sciences, 8(1):644–681, 2015. 3
- [34] Z. Peng, T. Wu, Y. Xu, M. Yan, and W. Yin. Coordinate friendly structures, algorithms and applications. *Annals of Mathematical Sciences and Applications*, 1(1):57–119, 2016. 3

- [35] Z. Peng, Y. Xu, M. Yan, and W. Yin. Arock: an algorithmic framework for asynchronous parallel coordinate updates. *SIAM Journal on Scientific Computing*, 38(5):A2851–A2879, 2016. 3
- [36] J.-C. Pesquet and A. Repetti. A class of randomized primal-dual algorithms for distributed optimization. arXiv preprint arXiv:1406.6404, 2014. 3
- [37] P. Richtárik and M. Takáč. Parallel coordinate descent methods for big data optimization.

 Mathematical Programming, pages 1–52, 2012. 2
- [38] P. Richtárik and M. Takáč. Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function. *Mathematical Programming*, 144(1-2):1–38, 2014. 2
- [39] Y. Xu. Hybrid jacobian and gauss-seidel proximal block coordinate update methods for linearly constrained convex programming. arXiv preprint arXiv:1608.03928, 2016. 3
- [40] Y. Xu. Accelerated first-order primal-dual proximal methods for linearly constrained composite convex programming. SIAM Journal on Optimization, 27(3):1459–1484, 2017. 3
- [41] Y. Xu. Asynchronous parallel primal-dual block update methods. arXiv preprint arXiv:1705.06391, 2017. 3

A Technical proofs: Section 2

In this section, we give the detailed proofs of the lemmas and theorems in section 2. The following lemma will be used a few times. Note that when S = [M], the result is deterministic.

Lemma A.1 Let S be a uniformly selected subset of [M] with cardinality m and x^o be a vector independent of S. Suppose x^+ is a random vector dependent on S and its coordinates out of S are the same as x^o . Let $\beta \in \mathbb{R}$, λ^o and r^o be vectors independent of S, and W a positive semidefinite $M \times M$ block diagonal matrix. If

$$\nabla_S f(x^o) + \tilde{\nabla} g_S(x_S^+) - A_S^\top (\lambda^o - \beta r^o) + W_S(x_S^+ - x_S^o) = 0,$$

then for any x, it holds that

$$\mathbb{E}_{S}\left[F(x^{+}) - F(x) + \frac{\mu}{2}\|x^{+} - x\|^{2} - \left\langle A(x^{+} - x), \lambda^{o} - \beta r^{o} \right\rangle\right]
\leq (1 - \theta) \left[F(x^{o}) - F(x) + \frac{\mu}{2}\|x^{o} - x\|^{2} - \left\langle A(x^{o} - x), \lambda^{o} - \beta r^{o} \right\rangle\right]
- \frac{1}{2}\mathbb{E}_{S}\left[\|x^{+} - x\|_{W}^{2} - \|x^{o} - x\|_{W}^{2} + \|x^{+} - x^{o}\|_{W - L_{m}I}^{2}\right],$$
(45)

where $\theta = \frac{m}{M}$, L_m is given in Assumption 4, and the expectation is taken on S.

Proof. For any x, we have

$$\left\langle x_S^+ - x_S, \nabla_S f(x^o) + \tilde{\nabla} g_S(x_S^+) - A_S^\top (\lambda^o - \beta r^o) + W_S(x_S^+ - x_S^o) \right\rangle = 0.$$

We split the left hand side of the above equation into four terms and bound each of them as below. First, we have

$$\mathbb{E}_{S} \left\langle x_{S}^{+} - x_{S}, \nabla_{S} f(x^{o}) \right\rangle
= \mathbb{E}_{S} \left\langle x^{+} - x^{o}, \nabla_{f}(x^{o}) \right\rangle + \mathbb{E}_{S} \left\langle x_{S}^{o} - x_{S}, \nabla_{S} f(x^{o}) \right\rangle
\geq \mathbb{E}_{S} \left[f(x^{+}) - f(x^{o}) - \frac{L_{m}}{2} \|x^{+} - x^{o}\|^{2} \right] + \theta[f(x^{o}) - f(x)]
= \mathbb{E}_{S} \left[f(x^{+}) - f(x) - \frac{L_{m}}{2} \|x^{+} - x^{o}\|^{2} \right] - (1 - \theta)[f(x^{o}) - f(x)], \tag{46}$$

where the first equality uses the fact $x_i^+ = x_i^o$, $\forall i \notin S$, and the inequality follows from the uniform distribution of S, the convexity of f, and also the inequality (18).

Secondly, it follows from the strong convexity of g that

$$\left\langle x_S^+ - x_S, \tilde{\nabla} g_S(x_S^+) \right\rangle \ge g_S(x_S^+) - g_S(x_S) + \sum_{i \in S} \frac{\mu}{2} \|x_i^+ - x_i\|^2.$$
 (47)

Since $g_S(x_S^+) - g_S(x_S) = g(x^+) - g(x^o) + g_S(x_S^o) - g_S(x_S)$ and $\mathbb{E}_S[g_S(x_S^o) - g_S(x_S)] = \theta[g(x^o) - g(x)]$, we have

$$\mathbb{E}_{S}[g_{S}(x_{S}^{+}) - g_{S}(x_{S})] = \mathbb{E}_{S}[g(x^{+}) - g(x^{o})] + \theta[g(x^{o}) - g(x)]$$

$$= \mathbb{E}_{S}[g(x^{+}) - g(x)] - (1 - \theta)[g(x^{o}) - g(x)]. \tag{48}$$

Similarly, it holds $\mathbb{E}_S \sum_{i \in S} \frac{\mu}{2} ||x_i^+ - x_i||^2 = \frac{\mu}{2} \left(\mathbb{E}_S ||x^+ - x||^2 - (1 - \theta) ||x^o - x||^2 \right)$. Hence, taking expectation on both sides of (47) yields

$$\mathbb{E}_{S} \left\langle x_{S}^{+} - x_{S}, \tilde{\nabla} g_{S}(x_{S}^{+}) \right\rangle$$

$$\geq \mathbb{E}_{S} \left[g(x^{+}) - g(x) + \frac{\mu}{2} \|x^{+} - x\|^{2} \right] - (1 - \theta) \left[g(x^{o}) - g(x) + \frac{\mu}{2} \|x^{o} - x\|^{2} \right]. \tag{49}$$

Thirdly, by essentially the same arguments on showing (48), we have

$$\mathbb{E}_S \left\langle x_S^+ - x_S, -A_S^\top (\lambda^o - \beta r^o) \right\rangle = -\mathbb{E}_S \left\langle A(x^+ - x), \lambda^o - \beta r^o \right\rangle + (1 - \theta) \left\langle A(x^o - x), \lambda^o - \beta r^o \right\rangle. \tag{50}$$

Fourth, note $\langle x_S^+ - x_S, W_S(x_S^+ - x_S^o) \rangle = \langle x^+ - x, W(x^+ - x^o) \rangle$, and thus by (6),

$$\mathbb{E}_{S}\left\langle x_{S}^{+} - x_{S}, W_{S}(x_{S}^{+} - x_{S}^{o})\right\rangle = \frac{1}{2}\mathbb{E}_{S}\left[\|x^{+} - x\|_{W}^{2} - \|x^{o} - x\|_{W}^{2} + \|x^{+} - x^{o}\|_{W}^{2}\right]. \tag{51}$$

The desired result is obtained by adding (46), (49), (50), and (51), and recalling F = f + g.

A.1 Proof of Lemma 2.1

From (7a), we have the optimality condition

$$\nabla f(x^k) - A^{\top}(\lambda^k - \beta_k r^k) + \tilde{\nabla} g(x^{k+1}) + P^k(x^{k+1} - x^k) = 0.$$

Hence, for any x such that Ax = b, it follows from the definition of Φ in (3) and Lemma A.1 with S = [M], $x^o = x^k$, $\lambda^o = \lambda^k$, $\beta = \beta_k$, $x^+ = x^{k+1}$, and $W = P^k$ that

$$\Phi(x^{k+1}, x, \lambda) \leq \left\langle Ax^{k+1} - b, \lambda^k - \beta_k r^k \right\rangle - \left\langle Ax^{k+1} - b, \lambda \right\rangle
- \frac{1}{2} \mathbb{E}_S \left[\|x^{k+1} - x\|_{P^k + \mu I}^2 - \|x^k - x\|_{P^k}^2 + \|x^{k+1} - x^k\|_{P^k - L_f I}^2 \right].$$
(52)

Using the fact $\lambda^{k+1} = \lambda^k - \rho_k(Ax^{k+1} - b)$, we have

$$\left\langle Ax^{k+1} - b, \lambda^k - \lambda \right\rangle = \frac{1}{\rho_k} \left\langle \lambda^k - \lambda^{k+1}, \lambda^k - \lambda \right\rangle$$

$$\stackrel{(6)}{=} \frac{1}{2\rho_k} \left[\|\lambda - \lambda^k\|^2 - \|\lambda - \lambda^{k+1}\|^2 + \|\lambda^k - \lambda^{k+1}\|^2 \right]. \tag{53}$$

In addition, we write $r^k = r^k - r^{k+1} + r^{k+1} = r^{k+1} - A(x^{k+1} - x^k)$ and have

$$\left\langle Ax^{k+1} - b, -\beta_k r^k \right\rangle
= -\beta_k \|r^{k+1}\|^2 + \beta_k \left\langle A(x^{k+1} - x), A(x^{k+1} - x^k) \right\rangle
\stackrel{(6)}{=} -\beta_k \|r^{k+1}\|^2 + \frac{\beta_k}{2} \left[\|A(x^{k+1} - x)\|^2 - \|A(x^k - x)\|^2 + \|A(x^{k+1} - x^k)\|^2 \right]$$
(54)

Substituting (53) and (54) into (52) gives the inequality in (8).

A.2 Proof of Theorem 2.2

First, we have

$$\sum_{k=1}^{t} \frac{k + k_0 + 1}{2\rho_k} \left[\|\lambda - \lambda^k\|^2 - \|\lambda - \lambda^{k+1}\|^2 \right] \\
= \frac{k_0 + 2}{2\rho_1} \|\lambda - \lambda^1\|^2 - \frac{t + k_0 + 1}{2\rho_t} \|\lambda - \lambda^{t+1}\|^2 + \sum_{k=2}^{t} \left(\frac{k + k_0 + 1}{2\rho_k} - \frac{k + k_0}{2\rho_{k-1}} \right) \|\lambda - \lambda^k\|^2 \\
\leq \frac{k_0 + 2}{2\rho_1} \|\lambda - \lambda^1\|^2.$$
(55)

In addition,

$$-\sum_{k=1}^{t} \frac{k+k_{0}+1}{2} \left(\|x^{k+1}-x\|_{P^{k}-\beta_{k}A^{\top}A+\mu I}^{2} - \|x^{k}-x\|_{P^{k}-\beta_{k}A^{\top}A}^{2} \right)$$

$$= \frac{k_{0}+2}{2} \|x^{1}-x\|_{P^{1}-\beta_{1}A^{\top}A}^{2} - \frac{t+k_{0}+1}{2} \|x^{t+1}-x\|_{P^{t}-\beta_{t}A^{\top}A+\mu I}^{2}$$

$$+\frac{1}{2} \sum_{k=2}^{t} \left((k+k_{0}+1) \|x^{k}-x\|_{P^{k}-\beta_{k}A^{\top}A}^{2} - (k+k_{0}) \|x^{k}-x\|_{P^{k-1}-\beta_{k-1}A^{\top}A+\mu I}^{2} \right)$$

$$\stackrel{(11)}{\leq} \frac{k_{0}+2}{2} \|x^{1}-x\|_{P^{1}-\beta_{1}A^{\top}A}^{2} - \frac{t+k_{0}+1}{2} \|x^{t+1}-x\|_{P^{t}-\beta_{t}A^{\top}A+\mu I}^{2}. \tag{56}$$

Now multiplying $k + k_0 + 1$ to both sides of (8) and adding it over k, we obtain (12) by using (55) and (56), and noting $\|\lambda^k - \lambda^{k+1}\|^2 = \rho_k^2 \|r^{k+1}\|^2$ and $\|x^{k+1} - x^k\|_{P^k - \beta_k A^\top A - L_f I}^2 \ge 0$.

A.3 Proof of Theorem 2.3

From the choice of k_0 and the condition $P - \beta A^{\top} A \leq \frac{\mu}{2} I$, it is not difficult to verify

$$(k + k_0 + 1) \left[kP - k\beta A^{\mathsf{T}} A + L_f I \right] \leq (k + k_0) \left[(k - 1)P - (k - 1)\beta A^{\mathsf{T}} A + (L_f + \mu)I \right], \forall k \geq 1.$$

Hence, the condition in (11) holds. In addition, it is easy to see that all conditions in (9) and (10) also hold. Therefore, we have (12), which, by taking parameters in (14) and $x = x^*$, reduces to

$$\sum_{k=1}^{t} (k + k_0 + 1) \Phi(x^{k+1}, x^*, \lambda) + \sum_{k=1}^{t} \frac{k(k + k_0 + 1)}{2} \beta \|r^{k+1}\|^2 + \frac{t + k_0 + 1}{2} \|x^{t+1} - x^*\|_{t(P - \beta A^{\top} A) + (L_f + \mu)I}^2 \leq \phi_1(x^*, \lambda), \tag{57}$$

where we have used the fact $\lambda^1 = 0$.

Letting $\lambda = \lambda^*$, we have from (5) and (57) that (by dropping nonnegative $\Phi(x^{k+1}, x^*, \lambda^*)$'s):

$$\frac{t(t+k_0+1)}{2}\beta \|r^{t+1}\|^2 + \frac{t+k_0+1}{2} \|x^{t+1} - x^*\|_{t(P-\beta A^\top A) + (L_f + \mu)I}^2 \le \phi_1(x^*, \lambda^*),$$

which indicates (15). In addition, from the convexity of F and (57), we have that for any λ , it holds $\frac{t(t+2k_0+3)}{2}\Phi(\bar{x}^{t+1},x^*,\lambda) \leq \phi_1(x^*,\lambda)$, which together with Lemmas 1.2 and 1.3 implies (16).

B Technical proofs: Section 3

In this section, we give the proofs of the lemmas and theorems in section 3.

B.1 Proof of Lemma 3.1

From the update in (17a), we have the optimality condition:

$$\nabla_{S_k} f(x^k) - A_{S_k}^{\top} (\lambda^k - \beta_k r^k) + \tilde{\nabla} g_{S_k} (x_{S_k}^{k+1}) + \eta_k (x_{S_k}^{k+1} - x_{S_k}^k) = 0.$$
 (58)

It follows from the update rule of λ that

$$-\langle Ax^{k+1} - b, \lambda^k \rangle = -\langle Ax^{k+1} - b, \lambda^{k+1} \rangle - \rho_k ||r^{k+1}||^2$$

Plugging (54) and the above equation into (45) with $S = S_k$, $\lambda^o = \lambda^k$, $\beta = \beta_k$, $x^o = x^k$, $x^+ = x^{k+1}$, $W = \eta_k I$, and x satisfying Ax = b, we have the desired result by taking expectation and recalling the definition of Δ in (2) and Φ in (3).

B.2 Proof of Theorem 3.2

Let $\beta_k = \beta$, $\rho_k = \rho$ and $\eta_k = \eta$ in (19), and also note $\mu = 0$ and $\eta \ge L_m + \beta ||A||^2$. We have

$$\begin{split} & \mathbb{E}\left[\Phi(x^{k+1}, x, \lambda^{k+1}) + (\beta - \rho) \|r^{k+1}\|^2\right] \\ & \leq (1 - \theta) \mathbb{E}\left[\Phi(x^k, x, \lambda^k) + \beta \|r^k\|^2\right] - \frac{1}{2} \mathbb{E}\left[\|x^{k+1} - x\|_{\eta I - \beta A^\top A}^2 - \|x^k - x\|_{\eta I - \beta A^\top A}^2\right]. \end{split}$$

Summing the above inequality over k=1 through t and noting $\rho \leq \theta \beta$ give

$$\mathbb{E}\left[\Phi(x^{t+1}, x, \lambda^{t+1}) + (\beta - \rho) \|r^{t+1}\|^{2}\right] + \theta \sum_{k=1}^{t-1} \mathbb{E}\Phi(x^{k+1}, x, \lambda^{k+1})
\leq (1 - \theta) \mathbb{E}\left[\Phi(x^{1}, x, \lambda^{1}) + \beta \|r^{1}\|^{2}\right] + \frac{1}{2} \|x^{1} - x\|_{\eta I - \beta A^{\top} A}^{2}.$$
(59)

By the update of λ , it follows that

$$\theta \sum_{k=1}^{t-1} \Phi(x^{k+1}, x, \lambda^{k+1}) = \theta \sum_{k=1}^{t-1} \left[\Phi(x^{k+1}, x, \lambda) + \frac{1}{\rho} \langle \lambda^{k+1} - \lambda, \lambda^{k+1} - \lambda^k \rangle \right]$$

$$= \theta \sum_{k=1}^{t-1} \Phi(x^{k+1}, x, \lambda) + \frac{\theta}{2\rho} \sum_{k=1}^{t-1} \left[\|\lambda^{k+1} - \lambda\|^2 - \|\lambda^k - \lambda\|^2 + \|\lambda^{k+1} - \lambda^k\|^2 \right]$$

$$= \theta \sum_{k=1}^{t-1} \Phi(x^{k+1}, x, \lambda) + \frac{\theta}{2\rho} \left[\|\lambda^t - \lambda\|^2 - \lambda^1 - \lambda\|^2 + \sum_{k=1}^{t-1} \|\lambda^{k+1} - \lambda^k\|^2 \right]$$
(60)

and

$$\Phi(x^{t+1}, x, \lambda^{t+1}) = \Phi(x^{t+1}, x, \lambda) - \langle \lambda^t - \lambda - \rho r^{t+1}, r^{t+1} \rangle
= \Phi(x^{t+1}, x, \lambda) - \langle \lambda^t - \lambda, r^{t+1} \rangle + \rho \|r^{t+1}\|^2.$$
(61)

Since $\rho \leq \theta \beta$, by Young's inequality, it holds

$$\beta \|r^{t+1}\|^2 - \langle \lambda^t - \lambda, r^{t+1} \rangle + \frac{\theta}{2\rho} \|\lambda^t - \lambda\|^2 \ge 0.$$

Then plugging (60) and (61) into (59), we have

$$\mathbb{E}\Phi(x^{t+1}, x, \lambda) + \theta \sum_{k=1}^{t-1} \mathbb{E}\Phi(x^{k+1}, x, \lambda)$$

$$\leq (1 - \theta) \mathbb{E}\left[\Phi(x^{1}, x, \lambda^{1}) + \beta \|r^{1}\|^{2}\right] + \frac{1}{2} \|x^{1} - x\|_{\eta I - \beta A^{\top} A}^{2} + \frac{\theta}{2\rho} \mathbb{E}\|\lambda^{1} - \lambda\|^{2}$$

$$\leq \mathbb{E}\phi_{2}(x, \lambda), \tag{62}$$

where in the last inequality we have used $\lambda^1 = 0$, $\theta > 0$ and $||r^1||^2 = ||x^1 - x||^2_{\beta A^\top A}$.

Therefore, from the convexity of F, it follows that $\mathbb{E}\Phi(\bar{x}^t, x^*, \lambda) \leq \frac{1}{1+\theta(t-1)}\mathbb{E}\phi_2(x^*, \lambda)$, $\forall \lambda$, and we obtain the desired result from Lemmas 1.2 and 1.3.

B.3 Proof of Theorem 3.3

We first establish a few inequalities below.

Proposition B.1 If (21e), (21f) and (21g) hold, then

$$-\sum_{k=1}^{t} (k+k_{0}+1) \mathbb{E}\left[\Delta_{\eta_{k}I-\beta_{k}A^{T}A}(x^{k+1}, x^{k}, x) - \frac{L_{m}}{2} \|x^{k+1} - x^{k}\|^{2}\right]$$

$$-\frac{\mu(t+k_{0}+1)}{2} \mathbb{E}\|x^{t+1} - x\|^{2} - \sum_{k=2}^{t} \frac{\mu(\theta(k+k_{0}+1)-1)}{2} \mathbb{E}\|x^{k} - x\|^{2}$$

$$\leq \frac{\eta_{1}(k_{0}+2)}{2} \mathbb{E}\|x^{1} - x\|^{2} - \frac{(t+k_{0}+1)}{2} \mathbb{E}\|x^{t+1} - x\|_{(\mu+\eta_{t})I-\beta_{t}A^{T}A}^{2}.$$

$$(63)$$

Proof. This inequality can be easily shown by noting that for any $1 \le k \le t$, the weight matrix of $\frac{1}{2} \|x^{k+1} - x^k\|^2$ is $\beta_k(k+k_0+1)A^{\top}A - (k+k_0+1)(\eta_k - L_m)I$, which is negative semidefinite, and for any $2 \le k \le t$, the weight matrix of $\frac{1}{2} \|x^k - x\|^2$ is

$$\left[\beta_{k-1}(k+k_0) - \beta_k(k+k_0+1)\right] A^{\top} A + \left[(k+k_0+1)\eta_k - (k+k_0)\eta_{k-1} - \mu(\theta(k+k_0+1)-1)\right] I,$$
 which is also negative semidefinite.

Proposition B.2 If (21a), (21c) and (21d) hold, then

$$-\frac{t+k_{0}+1}{\rho_{t}}\mathbb{E}\Delta(\lambda^{t+1},\lambda^{t},\lambda) - \sum_{k=2}^{t} \frac{\theta(k+k_{0}+1)-1}{\rho_{k-1}}\mathbb{E}\Delta(\lambda^{k},\lambda^{k-1},\lambda)$$

$$\leq \frac{\theta(k_{0}+3)-1}{2\rho_{1}}\mathbb{E}\|\lambda^{1}-\lambda\|^{2}.$$
(64)

Proof. On the left hand side of (64), the coefficient of each $\frac{1}{2}\|\lambda^{k+1} - \lambda^k\|^2$ is negative. For $2 \le k \le t-1$, the coefficient of $\frac{1}{2}\|\lambda^k - \lambda\|^2$ is $\frac{\theta(k+k_0+2)-1}{\rho_k} - \frac{\theta(k+k_0+1)-1}{\rho_{k-1}}$, which is nonpositive; the coefficient of $\frac{1}{2}\|\lambda^t - \lambda\|^2$ is $\frac{t+k_0+1}{\rho_t} - \frac{\theta(t+k_0+1)-1}{\rho_{t-1}}$, which is nonpositive; the coefficient of $\frac{1}{2}\|\lambda^{t+1} - \lambda\|^2$ is also nonpositive. Hence, dropping these nonpositive terms, we have the desired result.

Now we are ready to prove Theorem 3.3.

Proof. [of Theorem 3.3]

Multiplying $k + k_0 + 1$ to both sides of (19), summing it up from k = 1 through t, and moving the terms about $\Phi(x^k, x, \lambda^k) + \frac{\mu}{2} ||x^k - x||^2$ and $||r^k||^2$ to the left hand side for $2 \le k \le t$ give

$$(t+k_{0}+1)\mathbb{E}\left[\Phi(x^{t+1},x,\lambda^{t+1})+(\beta_{t}-\rho_{t})\|r^{t+1}\|^{2}+\frac{\mu}{2}\|x^{t+1}-x\|^{2}\right]$$

$$+\sum_{k=2}^{t}\left(\theta(k+k_{0}+1)-1\right)\mathbb{E}\left[\Phi(x^{k},x,\lambda^{k})+\frac{\mu}{2}\|x^{k}-x\|^{2}\right]$$

$$+\sum_{k=2}^{t}\left((\beta_{k-1}-\rho_{k-1})(k+k_{0})-(1-\theta)(k+k_{0}+1)\beta_{k}\right)\mathbb{E}\|r^{k}\|^{2}$$

$$\leq (1-\theta)(k_{0}+2)\mathbb{E}\left[\Phi(x^{1},x,\lambda^{1})+\beta_{1}\|r^{1}\|^{2}+\frac{\mu}{2}\|x^{1}-x\|^{2}\right]$$

$$-\sum_{k=1}^{t}(k+k_{0}+1)\mathbb{E}\left[\Delta_{\eta_{k}I-\beta_{k}A^{T}A}(x^{k+1},x^{k},x)-\frac{L_{m}}{2}\|x^{k+1}-x^{k}\|^{2}\right].$$
(65)

Hence, from (21b) and (63), it follows that

$$(t+k_{0}+1)\mathbb{E}\Phi(x^{t+1},x,\lambda^{t+1}) + \sum_{k=2}^{t} (\theta(k+k_{0}+1)-1)\mathbb{E}\Phi(x^{k},x,\lambda^{k})$$

$$\leq (1-\theta)(k_{0}+2)\mathbb{E}\left[\Phi(x^{1},x,\lambda^{1}) + \beta_{1}\|r^{1}\|^{2} + \frac{\mu}{2}\|x^{1}-x\|^{2}\right]$$

$$+ \frac{\eta_{1}(k_{0}+2)}{2}\mathbb{E}\|x^{1}-x\|^{2} - \frac{t+k_{0}+1}{2}\mathbb{E}\|x^{t+1}-x\|_{(\mu+\eta_{t})I-\beta_{t}A^{T}A}^{2}.$$

$$(66)$$

In addition, from the update of λ in (17c), we have

$$\langle \lambda^{k+1} - \lambda, Ax^{k+1} - b \rangle = -\frac{1}{\rho_k} \langle \lambda^{k+1} - \lambda, \lambda^{k+1} - \lambda^k \rangle = -\frac{1}{\rho_k} \Delta(\lambda^{k+1}, \lambda^k, \lambda), \tag{67}$$

and thus

$$(t + k_0 + 1)\mathbb{E}\langle\lambda^{t+1} - \lambda, Ax^{t+1} - b\rangle + \sum_{k=2}^{t} (\theta(k + k_0 + 1) - 1)\mathbb{E}\langle\lambda^k - \lambda, Ax^k - b\rangle$$

$$= -\frac{t + k_0 + 1}{\rho_t} \mathbb{E}\Delta(\lambda^{t+1}, \lambda^t, \lambda) - \sum_{k=2}^{t} \frac{\theta(k + k_0 + 1) - 1}{\rho_{k-1}} \mathbb{E}\Delta(\lambda^k, \lambda^{k-1}, \lambda)$$

$$\stackrel{(64)}{\leq} \frac{\theta(k_0 + 3) - 1}{2\rho_1} \mathbb{E}\|\lambda^1 - \lambda\|^2.$$

Since $\Phi(x^k, x, \lambda) = \Phi(x^k, x, \lambda^k) + \langle \lambda^k - \lambda, Ax^k - b \rangle$, we obtain the desired result by adding the above inequality to (66).

B.4 Proof of Proposition 3.4

Note that (24) implies $k_0 \ge \frac{4}{\theta}$, and thus (21a) must hold. Also, it is easy to see that (21d) holds with equality from the second equation of (23b). Since $I \succeq \frac{A^{\top}A}{\|A\|_2^2}$, we can easily have (21f) by plugging in β_k and η_k defined in (23a) and (23c) respectively.

To verify (21c), we plug in ρ_k defined in the first equation of (23b), and it is equivalent to requiring that for any $2 \le k \le t - 1$

$$\frac{\theta(k+k_0+1)-1}{\theta(k-1)+2+\theta} \ge \frac{\theta(k+k_0+2)-1}{\theta(k+2+\theta)} \iff 1 + \frac{\theta(k_0+1)-3}{\theta(k+2)} \ge 1 + \frac{\theta(k_0+1)-3}{\theta(k+2+\theta)}.$$

The inequality on the right hand side obviously holds, and thus we have (21c).

Plugging in the formula of β_k , (21e) is equivalent to

$$(\theta k + 2 + \theta)(k + k_0 + 1) \ge (\theta k + 2)(k + k_0),$$

which holds trivially, and thus (21e) follows.

With the given β_k and ρ_k , (21b) becomes $\frac{6}{6-5\theta}(\theta k+2)(k+k_0) \geq (k+k_0+1)(\theta k+2+\theta)$, $\forall 2 \leq k \leq t$, which is equivalent to $\frac{6}{6-5\theta} \geq \frac{(k_0+3)(3\theta+2)}{(k_0+2)(2\theta+2)}$. Note that $\frac{k_0+3}{k_0+2}$ is decreasing with respect to $k_0 \geq 0$ and also $\frac{6}{6-5\theta} \geq \frac{(\frac{3}{\theta}+3)(3\theta+2)}{(\frac{3}{\theta}+2)(2\theta+2)}$. Hence, (21b) is satisfied from the fact $k_0 \geq \frac{4}{\theta}$.

Finally, we show (21g). Plugging in η_k , we have that (21g) becomes

$$(k+k_0)\left(\frac{\mu}{2}(\theta k+2)+L_m\right)+\mu(\theta(k+k_0+1)-1)\geq (k+k_0+1)\left(\frac{\mu}{2}(\theta k+2+\theta)+L_m\right), \forall k\geq 2,$$

which is equivalent to $k_0 + 1 \ge \frac{4}{\theta} + \frac{2L_m}{\theta\mu}$. Hence, for k_0 given in (24), (21g) must hold. Therefore, we have verified all conditions in (21).

B.5 Proof of Theorem 3.5

From Proposition 3.4, we have the inequality in (22) that, as $\lambda^1 = 0$, reduces to

$$(t+k_0+1)\mathbb{E}\Phi(x^{t+1},x,\lambda) + \sum_{k=2}^{t} (\theta(k+k_0+1)-1)\mathbb{E}\Phi(x^k,x,\lambda)$$

$$\leq \phi_3(x,\lambda) - \frac{t+k_0+1}{2}\mathbb{E}\|x^{t+1}-x\|_{(\mu+\eta_t)I-\beta_tA^\top A}^2.$$
(68)

For $\rho \geq 1$, we have

$$(\mu + \eta_t)I - \beta_t A^{\top} A \succeq \left(\frac{(\rho - 1)\mu}{2\rho}(\theta t + \theta + 2) + \mu + L_m\right)I. \tag{69}$$

Letting $x = x^*$ and using the convexity of F, we have from (68) and the above inequality that

$$\mathbb{E}\left[F(\bar{x}^{t+1}) - F(x^*) - \left\langle \lambda, A\bar{x}^{t+1} - b\right\rangle\right] \le \frac{1}{T} \mathbb{E}\phi_3(x^*, \lambda), \,\forall \lambda,\tag{70}$$

which together with Lemmas 1.2 and 1.3 with $\gamma = \max(2\|\lambda^*\|, 1 + \|\lambda^*\|)$ indicates (25).

In addition, note

$$\Phi(x^{t+1}, x^*, \lambda^*) \ge \frac{\mu}{2} ||x^{t+1} - x^*||^2.$$

Hence, letting $(x, \lambda) = (x^*, \lambda^*)$ in (68) and using (5), we have from (69) that

$$\frac{t+k_0+1}{2} \left(\frac{(\rho-1)\mu}{2\rho} (\theta t + \theta + 2) + 2\mu + L_m \right) \mathbb{E} \|x^{t+1} - x^*\|^2 \le \phi_3(x^*, \lambda^*), \tag{71}$$

and the proof is completed.

C Technical proofs: Section 4

In this section, we provide the proofs of the lemmas and theorems in section 4.

C.1 Proof of Lemma 4.1

Note $r^{k+1} - r^k = A(x^{k+1} - x^k) + B(y^{k+1} - y^k)$. Hence by (6), we have

$$\left\langle A(x^{k+1} - x), -\beta r^k \right\rangle = -\beta \left\langle A(x^{k+1} - x), r^{k+1} \right\rangle + \beta \left\langle A(x^{k+1} - x), B(y^{k+1} - y^k) \right\rangle + \frac{\beta}{2} \left[\|A(x^{k+1} - x)\|^2 - \|A(x^k - x)\|^2 + \|A(x^{k+1} - x^k)\|^2 \right].$$
(72)

In addition, $\langle A(x^{k+1}-x), \lambda^k \rangle = \langle A(x^{k+1}-x), \lambda^{k+1} + \rho r^{k+1} \rangle$. Plugging this equation and (72) into (45) with $x^o = x^k, \lambda^o = \lambda^k, x^+ = x^{k+1}, W = \eta_x I$ and taking expectation yield

$$\mathbb{E}\left[F(x^{k+1}) - F(x) + \frac{\mu}{2} \|x^{k+1} - x\|^2 - \left\langle A(x^{k+1} - x), \lambda^{k+1} \right\rangle + (\beta - \rho) \left\langle A(x^{k+1} - x), r^{k+1} \right\rangle \right]
+ \frac{1}{2} \mathbb{E}\left[\|x^{k+1} - x\|_P^2 - \|x^k - x\|_P^2 + \|x^{k+1} - x^k\|_{P-L_m I}^2\right]
\leq (1 - \theta) \mathbb{E}\left[F(x^k) - F(x) + \frac{\mu}{2} \|x^k - x\|^2 - \left\langle A(x^k - x), \lambda^k - \beta r^k \right\rangle \right]
+ \beta \mathbb{E}\left\langle A(x^{k+1} - x), B(y^{k+1} - y^k) \right\rangle,$$
(73)

where $P = \eta_x I - \beta A^{\top} A$

From (30), the optimality condition for \tilde{y}^{k+1} is

$$\nabla h(\tilde{y}^{k+1}) - B^{\top} \lambda^k + \beta B^{\top} r^{k+\frac{1}{2}} + \eta_u(\tilde{y}^{k+1} - y^k) = 0.$$
 (74)

Since $\operatorname{Prob}(y^{k+1} = \tilde{y}^{k+1}) = \theta$, $\operatorname{Prob}(y^{k+1} = y^k) = 1 - \theta$, we have

$$\mathbb{E}\left\langle y^{k+1} - y, \nabla h(y^{k+1}) - B^{\top} \lambda^k + \beta B^{\top} r^{k+\frac{1}{2}} + \eta_y(y^{k+1} - y^k) \right\rangle$$
$$= (1 - \theta) \mathbb{E}\left\langle y^k - y, \nabla h(y^k) - B^{\top} \lambda^k + \beta B^{\top} r^{k+\frac{1}{2}} \right\rangle,$$

or equivalently,

$$\mathbb{E}\left\langle y^{k+1} - y, \nabla h(y^{k+1}) - B^{\top} \lambda^{k+1} + (\beta - \rho) B^{\top} r^{k+1} - \beta B^{\top} B(y^{k+1} - y^k) + \eta_y(y^{k+1} - y^k) \right\rangle
= (1 - \theta) \mathbb{E}\left\langle y^k - y, \nabla h(y^k) - B^{\top} \lambda^k + \beta B^{\top} r^k \right\rangle + \beta (1 - \theta) \mathbb{E}\left\langle B(y^k - y), A(x^{k+1} - x^k) \right\rangle. (75)$$

Recall $Q = \eta_y I - \beta B^{\top} B$. We have

$$\left\langle y^{k+1} - y, -\beta B^{\top} B(y^{k+1} - y^k) + \eta_y(y^{k+1} - y^k) \right\rangle = \frac{1}{2} \left[\|y^{k+1} - y\|_Q^2 - \|y^k - y\|_Q^2 + \|y^{k+1} - y^k\|_Q^2 \right].$$

Therefore adding (75) to (73), noting Ax + By = b, and plugging (67) with $\rho_k = \rho$, we have the desired result.

C.2 Proof of Theorem 4.2

Before proving Theorem 4.2, we establish a few inequalities. First, using Young's inequality, we have the following results.

Lemma C.1 For any $\tau_1, \tau_2 > 0$, it holds that

$$\langle A(x^{k+1} - x^*), B(y^{k+1} - y^k) \rangle \le \frac{1}{2\tau_1} ||A(x^{k+1} - x^*)||^2 + \frac{\tau_1}{2} ||B(y^{k+1} - y^k)||^2,$$
 (76)

$$\langle B(y^k - y^*), A(x^{k+1} - x^k) \rangle \le \frac{1}{2\tau_2} \|B(y^k - y^*)\|^2 + \frac{\tau_2}{2} \|A(x^{k+1} - x^k)\|^2.$$
 (77)

In addition, we are able to bound the λ -term by y-term and the residual r. The proofs are given in Appendix C.4 and C.5.

Lemma C.2 For any $\delta > 0$, we have

$$\mathbb{E}\|B^{\top}(\lambda^{k+1} - \lambda^{*})\|^{2} - (1 - \theta)(1 + \delta)\mathbb{E}\|B^{\top}(\lambda^{k} - \lambda^{*})\|^{2}
\leq 4\mathbb{E}\left[L_{h}^{2}\|y^{k+1} - y^{*}\|^{2} + \|Q(y^{k+1} - y^{k})\|^{2}\right] + 2(\beta - \rho)^{2}\mathbb{E}\|B^{\top}r^{k+1}\|^{2}
+2\rho^{2}(1 - \theta)(1 + \frac{1}{\delta})\mathbb{E}\left[\|B^{\top}r^{k+1}\|^{2} + \|B^{\top}B(y^{k+1} - y^{k})\|^{2}\right].$$
(78)

Lemma C.3 Assume (38). Then

$$\frac{\sigma_{\min}(BB^{\top})}{2} \left[\|\lambda^{k+1} - \lambda^*\|^2 - (1-\theta)\|\lambda^k - \lambda^*\|^2 + \frac{1}{\theta}\|\lambda^{k+1} - \lambda^k\|^2 \right]
\leq \|B^{\top}(\lambda^{k+1} - \lambda^*)\|^2 - (1-\theta)(1+\delta)\|B^{\top}(\lambda^k - \lambda^*)\|^2 + \kappa \|B^{\top}(\lambda^{k+1} - \lambda^k)\|^2, \tag{79}$$

where $\sigma_{\min}(BB^{\top})$ denotes the smallest singular value of BB^{\top} .

Lemma C.4 Let $c, \delta, \tau_1, \tau_2$ and κ be constants satisfying the conditions in Theorem 4.2. Then

$$\beta \mathbb{E} \langle A(x^{k+1} - x^*), B(y^{k+1} - y^k) \rangle + \beta (1 - \theta) \mathbb{E} \langle B(y^k - y^*), A(x^{k+1} - x^k) \rangle$$

$$+ \frac{c}{2} \sigma_{\min} (BB^{\top}) \mathbb{E} [\|\lambda^{k+1} - \lambda^*\|^2 - (1 - \theta)\|\lambda^k - \lambda^*\|^2 + \frac{1}{\theta} \|\lambda^{k+1} - \lambda^k\|^2]$$

$$\leq \frac{1}{2} \mathbb{E} \|x^{k+1} - x^k\|_{P-L_m I}^2 + \frac{\beta}{2\tau_1} \mathbb{E} \|A(x^{k+1} - x^*)\|^2$$

$$+ \frac{1}{2} \mathbb{E} \|y^{k+1} - y^k\|_Q^2 + \frac{\beta (1 - \theta)}{2\tau_2} \mathbb{E} \|B(y^k - y^*)\|^2 + 4cL_h^2 \mathbb{E} \|y^{k+1} - y^*\|^2$$

$$+ \left[c\rho^2 \left(\kappa + 2(1 - \theta)(1 + \frac{1}{\delta}) \right) + 2c(\beta - \rho)^2 \right] \mathbb{E} \|B^{\top} r^{k+1}\|^2.$$
(80)

Now we are ready to show Theorem 4.2.

Proof. [of Theorem 4.2]

Letting $(x, y, \lambda) = (x^*, y^*, \lambda^*)$ in (34), plugging (32) into it, and noting $Ax^* + By^* = b$, we have

$$\mathbb{E}\Psi(z^{k+1}, z^{*}) + (\beta - \rho)\mathbb{E}\|r^{k+1}\|^{2} + \mathbb{E}\left[\Delta_{P}(x^{k+1}, x^{k}, x^{*}) - \frac{L_{m}}{2}\|x^{k+1} - x^{k}\|^{2}\right] \\
+ \mathbb{E}\Delta_{Q}(y^{k+1}, y^{k}, y^{*}) + \frac{\mu}{2}\mathbb{E}\|x^{k+1} - x^{*}\|^{2} + \frac{1}{\rho}\mathbb{E}\Delta(\lambda^{k+1}, \lambda^{k}, \lambda^{*}) \\
\leq (1 - \theta)\mathbb{E}\Psi(z^{k}, z^{*}) + \beta(1 - \theta)\mathbb{E}\|r^{k}\|^{2} + \frac{1 - \theta}{\rho}\mathbb{E}\Delta(\lambda^{k}, \lambda^{k-1}, \lambda^{*}) + \frac{\mu(1 - \theta)}{2}\mathbb{E}\|x^{k} - x^{*}\|^{2} \\
+ \beta\mathbb{E}\langle A(x^{k+1} - x^{*}), B(y^{k+1} - y^{k})\rangle + \beta(1 - \theta)\mathbb{E}\langle B(y^{k} - y^{*}), A(x^{k+1} - x^{k})\rangle, \tag{81}$$

where Ψ is defined in (36). Note

$$\begin{split} &\frac{1}{\rho}\Delta(\lambda^{k+1},\lambda^k,\lambda^*)\\ &=\frac{1}{2\rho}\big[\|\lambda^{k+1}-\lambda^*\|^2-(1-\theta)\|\lambda^k-\lambda^*\|^2+\frac{1}{\theta}\|\lambda^{k+1}-\lambda^k\|^2\big]-\frac{\rho}{2}(\frac{1}{\theta}-1)\|r^{k+1}\|^2-\frac{\theta}{2\rho}\|\lambda^k-\lambda^*\|^2, \end{split}$$

and

$$\begin{split} &\frac{1-\theta}{\rho}\Delta(\lambda^k,\lambda^{k-1},\lambda^*)\\ &=\frac{1}{2\rho}\big[\|\lambda^k-\lambda^*\|^2-(1-\theta)\|\lambda^{k-1}-\lambda^*\|^2+\frac{1}{\theta}\|\lambda^k-\lambda^{k-1}\|^2\big]-\frac{\rho}{2}(\frac{1}{\theta}-(1-\theta))\|r^k\|^2-\frac{\theta}{2\rho}\|\lambda^k-\lambda^*\|^2. \end{split}$$

Adding (80) to (81) and plugging the above two equations yield

$$\mathbb{E}\Psi(z^{k+1},z^*) + (\beta - \rho)\mathbb{E}\|r^{k+1}\|^2 + \mathbb{E}\left[\Delta_P(x^{k+1},x^k,x^*) - \frac{L_m}{2}\|x^{k+1} - x^k\|^2\right] \\ + \mathbb{E}\Delta_Q(y^{k+1},y^k,y^*) + \frac{\mu}{2}\mathbb{E}\|x^{k+1} - x^*\|^2 - \frac{\rho}{2}(\frac{1}{\theta} - 1)\mathbb{E}\|r^{k+1}\|^2 - \frac{\theta}{2\rho}\mathbb{E}\|\lambda^k - \lambda^*\|^2 \\ + \left(\frac{1}{2\rho} + \frac{c}{2}\sigma_{\min}(BB^{\top})\right)\mathbb{E}\left[\|\lambda^{k+1} - \lambda^*\|^2 - (1-\theta)\|\lambda^k - \lambda^*\|^2 + \frac{1}{\theta}\|\lambda^{k+1} - \lambda^k\|^2\right] \\ \leq (1-\theta)\mathbb{E}\Psi(z^k,z^*) + \beta(1-\theta)\mathbb{E}\|r^k\|^2 - \frac{\rho}{2}(\frac{1}{\theta} - (1-\theta))\mathbb{E}\|r^k\|^2 - \frac{\theta}{2\rho}\mathbb{E}\|\lambda^k - \lambda^*\|^2 \\ + \frac{1}{2\rho}\mathbb{E}\left[\|\lambda^k - \lambda^*\|^2 - (1-\theta)\|\lambda^{k-1} - \lambda^*\|^2 + \frac{1}{\theta}\|\lambda^k - \lambda^{k-1}\|^2\right] \\ + \frac{\mu(1-\theta)}{2}\mathbb{E}\|x^k - x^*\|^2 + \frac{1}{2}\mathbb{E}\|x^{k+1} - x^k\|_{P-L_mI}^2 + \frac{\beta}{2\tau_1}\mathbb{E}\|A(x^{k+1} - x^*)\|^2 \\ + \frac{1}{2}\mathbb{E}\|y^{k+1} - y^k\|_Q^2 + \frac{\beta(1-\theta)}{2\tau_2}\mathbb{E}\|B(y^k - y^*)\|^2 + 4cL_h^2\mathbb{E}\|y^{k+1} - y^*\|^2 \\ + \left[c\rho^2\left(\kappa + 2(1-\theta)(1 + \frac{1}{\delta})\right) + 2c(\beta - \rho)^2\right]\mathbb{E}\|B^{\top}r^{k+1}\|^2.$$

Using the definition in (2) to expand $\Delta_P(x^{k+1}, x^k, x^*)$ and $\Delta_Q(y^{k+1}, y^k, y^*)$ in the above inequality, and then rearranging terms, we have

$$\mathbb{E}\Psi(z^{k+1}, z^{*}) + \left((\beta - \rho) - \frac{\rho}{2} (\frac{1}{\theta} - 1) \right) \mathbb{E} \| r^{k+1} \|^{2} \\
- \left[c\rho^{2} \left(\kappa + 2(1 - \theta) (1 + \frac{1}{\delta}) \right) + 2c(\beta - \rho)^{2} \right] \mathbb{E} \| B^{\top} r^{k+1} \|^{2} \\
+ \mathbb{E} \left[\frac{1}{2} \| x^{k+1} - x^{*} \|_{P}^{2} + \frac{\mu}{2} \| x^{k+1} - x^{*} \|^{2} - \frac{\beta}{2\tau_{1}} \| A(x^{k+1} - x^{*}) \|^{2} \right] \\
+ \mathbb{E} \left[\frac{1}{2} \| y^{k+1} - y^{*} \|_{Q}^{2} - 4cL_{h}^{2} \| y^{k+1} - y^{*} \|^{2} \right] \\
+ \left(\frac{1}{2\rho} + \frac{c}{2} \sigma_{\min}(BB^{\top}) \right) \mathbb{E} \left[\| \lambda^{k+1} - \lambda^{*} \|^{2} - (1 - \theta) \| \lambda^{k} - \lambda^{*} \|^{2} + \frac{1}{\theta} \| \lambda^{k+1} - \lambda^{k} \|^{2} \right] \\
\leq (1 - \theta) \mathbb{E} \Psi(z^{k}, z^{*}) + \beta(1 - \theta) \mathbb{E} \| r^{k} \|^{2} - \frac{\rho}{2} (\frac{1}{\theta} - (1 - \theta)) \mathbb{E} \| r^{k} \|^{2} + \frac{1}{2} \mathbb{E} \| x^{k} - x^{*} \|_{P}^{2} \\
+ \frac{\mu(1 - \theta)}{2} \mathbb{E} \| x^{k} - x^{*} \|^{2} + \frac{1}{2} \mathbb{E} \| y^{k} - y^{*} \|_{Q}^{2} + \frac{\beta(1 - \theta)}{2\tau_{2}} \mathbb{E} \| B(y^{k} - y^{*}) \|^{2} \\
+ \frac{1}{2\rho} \mathbb{E} \left[\| \lambda^{k} - \lambda^{*} \|^{2} - (1 - \theta) \| \lambda^{k-1} - \lambda^{*} \|^{2} + \frac{1}{\theta} \| \lambda^{k} - \lambda^{k-1} \|^{2} \right]. \tag{82}$$

Since $\rho = \theta \beta$, it holds

$$(\beta - \rho) - \frac{\rho}{2}(\frac{1}{\theta} - 1) = \frac{\beta - \rho}{2}, \quad \beta(1 - \theta) - \frac{\rho}{2}(\frac{1}{\theta} - (1 - \theta)) \le \frac{\beta(1 - \theta)}{2},$$

and thus the inequality (82) implies

$$\mathbb{E}\Psi(z^{k+1},z^*) + \frac{\beta - \rho}{2}\mathbb{E}\|r^{k+1}\|^2 - \left[c\rho^2\left(\kappa + 2(1-\theta)\left(1 + \frac{1}{\delta}\right)\right) + 2c(\beta - \rho)^2\right]\mathbb{E}\|B^{\top}r^{k+1}\|^2$$

$$+\mathbb{E}\left[\frac{1}{2}\|x^{k+1} - x^*\|_P^2 + \frac{\mu}{2}\|x^{k+1} - x^*\|^2 - \frac{\beta}{2\tau_1}\|A(x^{k+1} - x^*)\|^2\right] \\
+\mathbb{E}\left[\frac{1}{2}\|y^{k+1} - y^*\|_Q^2 - 4cL_h^2\|y^{k+1} - y^*\|^2\right] \\
+\left(\frac{1}{2\rho} + \frac{c}{2}\sigma_{\min}(BB^\top)\right)\mathbb{E}\left[\|\lambda^{k+1} - \lambda^*\|^2 - (1-\theta)\|\lambda^k - \lambda^*\|^2 + \frac{1}{\theta}\|\lambda^{k+1} - \lambda^k\|^2\right] \\
\leq \psi(z^k, z^*; P, Q, \beta, \rho, c, \tau_2), \tag{83}$$

where ψ is defined in (37).

From (33), it follows that

$$(1 - \alpha)\Psi(z^{k+1}, z^*) + \frac{\alpha\mu}{2} \|x^{k+1} - x^*\|^2 + \alpha\nu \|y^{k+1} - y^*\|^2 \le \Psi(z^{k+1}, z^*).$$
 (84)

In addition, note that

$$\begin{split} \|r^{k+1}\|^2 &= \|Ax^{k+1} + By^{k+1} - (Ax^* + By^*)\|^2 \\ &\leq 2\|A\|_2^2 \|x^{k+1} - x^*\|^2 + 2\|B\|_2^2 \|y^{k+1} - y^*\|^2 \\ &\leq \gamma \left(\frac{\alpha\mu}{4} \|x^{k+1} - x^*\|^2 + \frac{\alpha\nu}{4} \|y^{k+1} - y^*\|^2\right), \end{split}$$

and thus

$$\frac{1}{\gamma} \|r^{k+1}\|^2 \le \frac{\alpha \mu}{4} \|x^{k+1} - x^*\|^2 + \frac{\alpha \nu}{4} \|y^{k+1} - y^*\|^2.$$
 (85)

Adding (84) and (85) to (83) gives the desired result.

C.3 Proof of Theorem 4.3

From $0 < \alpha < \theta$, the full row-rankness of B, and the conditions in (41), it is easy to see that $\eta > 1$. Next we find lower bounds of the terms on the left hand of (40). Since $\eta \leq \frac{1-\alpha}{1-\theta}$, we have

$$\eta(1-\theta)\Psi(z^{k+1}, z^*) \le (1-\alpha)\Psi(z^{k+1}, z^*). \tag{86}$$

Note $||A||_2 \leq 1$ and

$$\left(\frac{\alpha\mu}{2} + \mu - \frac{\beta}{\tau_1}\right)I \succeq \frac{\frac{\alpha\mu}{2} + \theta\mu - \frac{\beta}{\tau_1}}{\eta_x + \mu(1-\theta)}(\eta_x I - \beta A^\top A) + \frac{\frac{\alpha\mu}{2} + \theta\mu - \frac{\beta}{\tau_1}}{\eta_x + \mu(1-\theta)}\mu(1-\theta)I + \mu(1-\theta)I.$$

Hence, from $\eta \leq 1 + \frac{\frac{\alpha\mu}{2} + \theta\mu - \frac{\beta}{\tau_1}}{\eta_x + \mu(1-\theta)}$ and $P = \eta_x I - \beta A^{\top} A$, it follows that

$$\eta \|x^{k+1} - x^*\|_{P+\mu(1-\theta)I}^2 \le \|x^{k+1} - x^*\|_{P+(\frac{\alpha\mu}{2} + \mu)I - \frac{\beta}{\tau_1}A^\top A}^2.$$
(87)

Similarly, since

$$\left(\frac{3\alpha\nu}{2} - 8cL_h^2\right)I \succeq \frac{\frac{3\alpha\nu}{2} - 8cL_h^2 - \frac{\beta(1-\theta)}{\tau_2}}{\eta_y + \frac{\beta(1-\theta)}{\tau_2}}(\eta_y I - \beta B^{\top}B) + \frac{\frac{3\alpha\nu}{2} - 8cL_h^2 - \frac{\beta(1-\theta)}{\tau_2}}{\eta_y + \frac{\beta(1-\theta)}{\tau_2}}\frac{\beta(1-\theta)}{\tau_2}I + \frac{\beta(1-\theta)}{\tau_2}I,$$

 $Q = \eta_u I - \beta B^{\top} B$, and $B^{\top} B \leq I$, we have

$$\eta \|y^{k+1} - y^*\|_{Q + \frac{\beta(1-\theta)}{\tau_2}B^\top B}^2 \le \|y^{k+1} - y^*\|_{Q + (\frac{3\alpha\nu}{2} - 8cL_h^2)I}^2. \tag{88}$$

For the r-term, we note from the definition of η that

$$\eta \frac{\beta(1-\theta)}{2} \le \left(\frac{\beta(1-\theta)}{2} + \frac{1}{\gamma}\right) - \left(c\rho^2\left(\kappa + 2(1-\theta)(1+\frac{1}{\delta})\right) + 2c(\beta-\rho)^2\right)$$

In addition, since $||B||_2 \leq 1$, it holds $||B^{\top}r^{k+1}|| \leq ||r^{k+1}||$, and thus

$$\eta \frac{\beta(1-\theta)}{2} \|r^{k+1}\|^2 \le \left(\frac{\beta(1-\theta)}{2} + \frac{1}{\gamma}\right) \|r^{k+1}\|^2 - \left(c\rho^2\left(\kappa + 2(1-\theta)(1+\frac{1}{\delta})\right) + 2c(\beta-\rho)^2\right) \|B^{\top}r^{k+1}\|^2. \tag{89}$$

Finally, it is obvious to have

$$\frac{\eta}{2\rho} \left[\|\lambda^{k+1} - \lambda^*\|^2 - (1-\theta)\|\lambda^k - \lambda^*\|^2 + \frac{1}{\theta}\|\lambda^{k+1} - \lambda^k\|^2 \right]
\leq \left(\frac{1}{2\rho} + \frac{c}{2} \sigma_{\min}(BB^\top) \right) \left[\|\lambda^{k+1} - \lambda^*\|^2 - (1-\theta)\|\lambda^k - \lambda^*\|^2 + \frac{1}{\theta}\|\lambda^{k+1} - \lambda^k\|^2 \right].$$
(90)

Therefore, we obtain (42) by the definition of ψ and adding (86) through (90).

C.4 Proof of Lemma C.2

Let $\tilde{\lambda}^{k+1} = \lambda^k - \rho(Ax^{k+1} + B\tilde{y}^{k+1} - b)$. Then from the update of y, we have

$$\mathbb{E}\|B^{\top}(\lambda^{k+1} - \lambda^*)\|^2 = \theta \mathbb{E}\|B^{\top}(\tilde{\lambda}^{k+1} - \lambda^*)\|^2 + (1 - \theta)\mathbb{E}\|B^{\top}(\lambda^k - \lambda^* - \rho(Ax^{k+1} + By^k - b))\|^2.$$
(91)

Below we bound the two terms on the right hand side of (91). First, the definition of $\tilde{\lambda}^{k+1}$ together with (74) implies

$$B^{\top} \tilde{\lambda}^{k+1} = \nabla h(\tilde{y}^{k+1}) + Q(\tilde{y}^{k+1} - y^k) + (\beta - \rho)B^{\top} (Ax^{k+1} + B\tilde{y}^{k+1} - b). \tag{92}$$

Hence, by the Young's inequality and the condition in (32b), we have

$$\theta \mathbb{E} \|B^{\top}(\tilde{\lambda}^{k+1} - \lambda^*)\|^2 \le 2\theta \mathbb{E} \|\nabla h(\tilde{y}^{k+1}) - \nabla h(y^*) + Q(\tilde{y}^{k+1} - y^k)\|^2 + 2\theta(\beta - \rho)^2 \mathbb{E} \|B^{\top}(Ax^{k+1} + B\tilde{y}^{k+1} - b)\|^2.$$
(93)

Since $\operatorname{Prob}(y^{k+1} = \tilde{y}^{k+1}) = \theta$ and $\operatorname{Prob}(y^{k+1} = y^k) = 1 - \theta$, it follows that

$$\mathbb{E}\|\nabla h(y^{k+1}) - \nabla h(y^*) + Q(y^{k+1} - y^k)\|^2$$

$$= \theta \mathbb{E}\|\nabla h(\tilde{y}^{k+1}) - \nabla h(y^*) + Q(\tilde{y}^{k+1} - y^k)\|^2 + (1 - \theta)\mathbb{E}\|\nabla h(y^k) - \nabla h(y^*)\|^2,$$

and thus

$$\theta \mathbb{E} \|\nabla h(\tilde{y}^{k+1}) - \nabla h(y^*) + Q(\tilde{y}^{k+1} - y^k)\|^2 \le \mathbb{E} \|\nabla h(y^{k+1}) - \nabla h(y^*) + Q(y^{k+1} - y^k)\|^2$$

Similarly,

$$\theta(\beta - \rho)^2 \mathbb{E} \|B^\top (Ax^{k+1} + B\tilde{y}^{k+1} - b)\|^2 \le (\beta - \rho)^2 \mathbb{E} \|B^\top (Ax^{k+1} + By^{k+1} - b)\|^2.$$

Plugging the above two equations into (93) and applying the Young's inequality and also the Lipschitz continuity of ∇h give

$$\theta \mathbb{E} \|B^{\top}(\tilde{\lambda}^{k+1} - \lambda^*)\|^2 \le 4\mathbb{E} \left[L_h^2 \|y^{k+1} - y^*\|^2 + \|Q(y^{k+1} - y^k)\|^2\right] + 2(\beta - \rho)^2 \mathbb{E} \|B^{\top} r^{k+1}\|^2. \tag{94}$$

In addition, from the Young's inequality, it follows for any $\delta > 0$ that

$$\|B^{\top}(\lambda^k - \lambda^* - \rho(Ax^{k+1} + By^k - b))\|^2 \le (1 + \delta)\|B^{\top}(\lambda^k - \lambda^*)\|^2 + \rho^2(1 + \frac{1}{\delta})\|B^{\top}(Ax^{k+1} + By^k - b)\|^2.$$

Note $||B^{\top}(Ax^{k+1} + By^k - b)||^2 \le 2||B^{\top}r^{k+1}||^2 + 2||B^{\top}B(y^{k+1} - y^k)||^2$. Therefore, plugging (94) and the above two inequalites into (91), we complete the proof.

C.5 Proof of Lemma C.3

It is straightforward to verify

$$||B^{\top}(\lambda^{k+1} - \lambda^*)||^2 - (1 - \theta)(1 + \delta)||B^{\top}(\lambda^k - \lambda^*)||^2 + \kappa ||B^{\top}(\lambda^{k+1} - \lambda^k)||^2$$

$$= \begin{bmatrix} \lambda^{k+1} - \lambda^* \\ \lambda^{k+1} - \lambda^k \end{bmatrix}^{\top} \begin{bmatrix} (1 - (1 - \theta)(1 + \delta)) & (1 - \theta)(1 + \delta) \\ (1 - \theta)(1 + \delta) & (\kappa - (1 - \theta)(1 + \delta)) \end{bmatrix} \otimes BB^{\top} \begin{bmatrix} (\lambda^{k+1} - \lambda^*) \\ (\lambda^{k+1} - \lambda^k) \end{bmatrix},$$

and

$$\begin{bmatrix} \lambda^{k+1} - \lambda^* \\ \lambda^{k+1} - \lambda^k \end{bmatrix}^{\top} \begin{bmatrix} \theta & (1-\theta) \\ (1-\theta) & (\frac{1}{\theta} - (1-\theta)) \end{bmatrix} \otimes I \begin{bmatrix} \lambda^{k+1} - \lambda^* \\ \lambda^{k+1} - \lambda^k \end{bmatrix}$$

$$= \begin{bmatrix} \|\lambda^{k+1} - \lambda^*\|^2 - (1-\theta)\|\lambda^k - \lambda^*\|^2 + \frac{1}{\theta}\|\lambda^{k+1} - \lambda^k\|^2 \end{bmatrix}.$$

Hence, we have the desired result from (38) and the inequality $U \otimes V \succeq \sigma_{\min}(V)U \otimes I$ for any PSD matrices U and V.

C.6 Proof of Lemma C.4

From (39a) and (39b), we have

$$\beta(1-\theta)\frac{\tau_2}{2}\|A(x^{k+1}-x^k)\|^2 \le \frac{1}{2}\|x^{k+1}-x^k\|_{P-L_mI}^2,$$

and

$$\begin{aligned} & 4c\|Q(y^{k+1}-y^k)\|^2 + 2c\rho^2(1-\theta)(1+\frac{1}{\delta})\|B^\top B(y^{k+1}-y^k)\|^2 + \frac{\beta\tau_1}{2}\|B(y^{k+1}-y^k)\|^2 \\ & \leq & \frac{1}{2}\|y^{k+1}-y^k\|_Q^2. \end{aligned}$$

The desired result is then obtained by adding the above two inequalities together with β times of (76), $\beta(1-\theta)$ times of (77), c times of both (78) and (79), and also noting $\lambda^{k+1} - \lambda^k = -\rho r^{k+1}$.