

# Randomized First-Order Methods for Saddle Point Optimization <sup>\*</sup>

Cong D. Dang <sup>†</sup>      Guanghui Lan <sup>‡</sup>

November 13, 2015

## Abstract

In this paper, we present novel randomized algorithms for solving saddle point problems whose dual feasible region is given by the direct product of many convex sets. Our algorithms can achieve an  $\mathcal{O}(1/N)$  and  $\mathcal{O}(1/N^2)$  rate of convergence, respectively, for general bilinear saddle point and smooth bilinear saddle point problems based on a new primal-dual termination criterion, and each iteration of these algorithms needs to solve only one randomly selected dual subproblem. Moreover, these algorithms do not require strongly convex assumptions on the objective function and/or the incorporation of a strongly convex perturbation term. They do not necessarily require the primal or dual feasible regions to be bounded or the estimation of the distance from the initial point to the set of optimal solutions to be available either. We show that when applied to linearly constrained problems, RPDs are equivalent to certain randomized variants of the alternating direction method of multipliers (ADMM), while a direct extension of ADMM does not necessarily converge when the number of blocks exceeds two.

**Keywords.** Stochastic Optimization, Block Coordinate Descent, Nonsmooth Optimization, Saddle Point Optimization, Alternating Direction Method of Multipliers

## 1 Introduction

Motivated by some recent applications in data analysis, there has been a growing interest in the design and analysis of randomized first-order methods for large-scale convex optimization. In these applications, the complex datasets are so big and often distributed over different storage locations. It is often impractical to assume that optimization algorithms can traverse an entire dataset once in each iteration, because doing so is either time consuming or unreliable, and often results in low resource utilization due to necessary synchronization among different computing units (e.g., CPUs, GPUs, and Cores) in a distributed computing environment. On the other hand, randomized algorithms can make progress by using information obtained from a randomly selected subset of data and thus provide much flexibility for their implementation in the aforementioned distributed environments.

---

<sup>\*</sup>This research was partially supported by NSF grants CMMI-1000347, CMMI-1254446, DMS-1319050, and ONR grant N00014-13-1-0036.

<sup>†</sup>Department of Industrial and Systems Engineering, University of Florida, Gainesville, FL, 32611. (email: [congdd@ufl.edu](mailto:congdd@ufl.edu)).

<sup>‡</sup>Department of Industrial and Systems Engineering, University of Florida, Gainesville, FL, 32611. (email: [glan@ise.ufl.edu](mailto:glan@ise.ufl.edu)).

In this paper, we focus on the development of randomized algorithms for solving a class of saddle point problems given by

$$\min_{x \in X} \left\{ h(x) + \max_{y \in Y} \langle Ax, y \rangle - J(y) \right\}, \quad (1.1)$$

where  $X \subseteq \mathbb{R}^n$  and  $Y \subseteq \mathbb{R}^m$  are closed convex sets,  $h : X \rightarrow \mathbb{R}$  and  $J : Y \rightarrow \mathbb{R}$  are closed convex functions, and  $A : \mathbb{R}^n \rightarrow \mathbb{R}^m$  denotes a given linear operator. Throughout this paper, we assume that

$$y = (y_1; \dots; y_p), \quad Y = Y_1 \times \dots \times Y_p, \quad \text{and} \quad J(y) = J_1(y_1) + \dots + J_p(y_p). \quad (1.2)$$

Here  $y_i \in Y_i$ ,  $i = 1, \dots, p$ ,  $Y_i \subseteq \mathbb{R}^{m_i}$  are given closed convex sets such that  $\sum_{i=1}^p m_i = m$ , and  $J_i : Y_i \rightarrow \mathbb{R}$ ,  $i = 1, \dots, p$ , are closed convex functions. Accordingly, we denote  $A = (A_1; \dots; A_p)$ , where  $A_i$  are given linear operators from  $\mathbb{R}^n$  to  $\mathbb{R}^{m_i}$ ,  $i = 1, \dots, p$ .

Problem (1.1)-(1.2) covers a few interesting subclasses of problems in the literature. One prominent example is to minimize the summation of several separable convex functions over some coupled linear constraints. Indeed, letting  $X = \mathbb{R}^n$  and  $h(x) = -b^T x$ , one can view problem (1.1)-(1.2) as the saddle-point reformation of

$$\begin{aligned} \min \quad & J_1(y_1) + J_2(y_2) + \dots + J_p(y_p) \\ \text{s.t.} \quad & A_1^T y_1 + A_2^T y_2 + \dots + A_p^T y_p = b, \\ & y_i \in Y_i, i = 1, \dots, p. \end{aligned} \quad (1.3)$$

The above problem has found wide applications in machine learning and image processing, and many first-order algorithms have developed for its solutions. More specifically, one can apply Nesterov's smoothing scheme [38], the primal-dual method [6, 9], and the mirror-prox method [36, 34, 8] to solve the saddle-point reformulation in (1.1). We can also apply some classic penalty-based approaches for solving (1.3). In particular, Lan and Monteiro discussed the complexity of first-order quadratic penalty methods [27] and augmented Lagrangian penalty methods [26] applied to problem (1.3). More recently, He, Juditsky and Nemirovski generalized the mirror-prox algorithm for solving problem (1.3) based on the exact penalty method [22]. When  $p = 2$ , a special augmented Lagrangian method, namely the alternating direction method of multipliers (ADMM) [11, 43, 29, 15, 12], has been intensively studied recently [3, 18, 19, 20, 41]. However, as shown by Chen et al. [7], a direction extension of ADMM does not necessarily converge when  $p > 2$ , unless some strong convexity assumptions on  $J_i$  and full row rank assumption on  $A_i$  are made (e.g., [20, 23, 48]). Observe that all these methods need to perform  $p$  projection subproblems over the sets  $Y_i$ ,  $i = 1, \dots, p$ , in every iteration.

Another interesting example is to minimize the regularized loss function given by

$$\min_{x \in X} h(x) + \sum_{i=1}^p f_i(A_i x), \quad (1.4)$$

where  $f_i : \mathbb{R}^{m_i} \rightarrow \mathbb{R}$  are closed convex functions with conjugate  $f_i^*$ ,  $i = 1, \dots, p$ . Clearly, problem (1.4) can be viewed as a special case of problem (1.1) with  $J_i = f_i^*$  and  $Y_i = \mathbb{R}^{m_i}$ ,  $i = 1, \dots, p$ . While the algorithms for solving problem (1.3) are mostly deterministic, much effort has been devoted to randomized first-order methods for solving problem (1.4), which can make progress by utilizing the (sub)gradient of a randomly selected component  $f_i(A_i x)$  only. More specifically, if  $f_i$  are general nonsmooth convex functions, one can apply the mirror-descent stochastic approximation in [37] or

the accelerated stochastic approximation in [24], which exhibit an  $\mathcal{O}(1/\sqrt{N})$  rate of convergence for solving problem (1.4). Here  $N$  denotes the number of iterations. Recently, some interesting development has been made [44, 2, 31, 47] under the assumption that  $f_i$  are smooth convex functions. Based on incremental averaging gradient method [2], Schmidt et. al. [44] developed a stochastic averaging gradient method and show that it exhibits an  $\mathcal{O}(1/N)$  rate of convergence for smooth problems and an linear rate of convergence for the case when  $f_i$  are smooth and strongly convex. This algorithm is also closely related to the stochastic dual coordinate ascent [46], a randomized version of dual coordinate ascent applied to the dual of problem (1.4) when  $h$  is strongly convex, see [39, 28, 40, 42, 1, 30, 10] for some recent developments on block coordinate descent methods.

In this paper, we propose a novel algorithm, namely the randomized primal-dual method, to solve problems in the form of (1.1)-(1.2). The main idea is to incorporate a block decomposition of dual space into the primal-dual algorithm in [6]. At each iteration, our algorithm requires to solve only one subproblem in dual space instead of  $p$  subproblems as in the primal-dual algorithm. By using a new primal-dual termination criterion inspired by the one employed by Monteiro and Svaiter [33], we show that our algorithm can achieve an  $\mathcal{O}(1/N)$  and  $\mathcal{O}(1/N^2)$  rate of convergence, respectively for solving general bilinear saddle point problems (without any strongly convex assumptions) and smooth bilinear saddle point problems (with  $J$  being strongly convex), where  $N$  is the number of iterations. Furthermore, we demonstrate that our algorithm can deal with the situation when either  $X$  or  $Y$  is unbounded, as long as a saddle point of problem (1.1)-(1.2) exists. It should be noted that these complexity results will have an extra constant factor which depends on the number of blocks  $p$ , but such a dependence is mild if  $p$  is not too big. In addition, we discuss possible extensions of the RPD method to the non-Euclidean geometry and also show that RPD applied to the linearly constrained problems in (1.3) is equivalent to a certain randomized variant of the ADMM method. To the best of our knowledge, all these developments seem to be new in the literature. In fact, our proof for the convergence of the ergodic mean of the primal-dual method for smooth bilinear saddle point problems was also new even under the deterministic setting (i.e.,  $p = 1$ ), <sup>1</sup>.

It should be noted that in a concurrent and independent work, Zhang and Xiao [49] presented a randomized version of the primal-dual method for solving a special class of regularized empirical risk minimization (ERM) problems given in the form of (1.4) <sup>2</sup>. However, the algorithms, analysis and termination criteria in these papers are significantly different: (a) our primal-dual algorithm does not involve any extrapolation step as used in [49]; (b) we employed a new primal-dual optimality gap to assess the quality of a feasible solution to problem (1.1), while [49] employs the distance to the optimal solution as the termination criterion; and (c) as a consequence, the convergence analyses in these papers are significantly different. In fact, the basic algorithm in [49] was designed for problems where  $h$  is strongly convex problems (similarly to those randomized dual coordinate descent methods [45]). Otherwise, one has to add a strongly convex perturbation to the objective function and impose stronger assumptions about  $f_i$  and  $h$ . Such a perturbation term can be properly chosen only if there exists a bound on the distance from the initial point to the set of optimal solutions, and hence are not best suitable for the linearly constrained problems in (1.3). In fact, the authors were not aware of the existence of any other randomized algorithms in the literature that do not require the incorporation of a perturbation term for solving (1.1)-(1.2), but can achieve the optimal rate of convergence in

<sup>1</sup>It is worth noting that Chambolle and Pock [?] had also released their results on the convergence of the ergodic means for deterministic primal-dual methods shortly after we released the initial version of the current paper in Sep., 2014.

<sup>2</sup>Note that [49] was also initially released in Sep., 2014.

terms of their dependence on  $N$  as shown in this paper.

This paper is organized as follows. We first discuss some new primal-dual termination criteria in Section 2. We then present a general RPD method in Section 3, and discuss its convergence properties for general bilinear saddle point and smooth bilinear saddle point problems under the assumption that the primal and dual feasible regions are bounded. In Section 3, we generalize the RPD method for the case when the feasible regions are unbounded and incorporate non-Euclidean distance generating functions into the RPD method. In Section 4, we discuss the relation of the RPD method to ADMM. Finally some brief concluding remarks are provided in Section 5.

## 2 The problem of interest and its termination criteria

We introduce in this section a few termination criteria that will be used to evaluate the solution quality for problem (1.1).

Denote  $Z \equiv X \times Y$ . For a given  $\hat{z} = (\hat{x}, \hat{y}) \in Z$ , let us define the gap function  $Q_0$  by

$$Q_0(\hat{z}, z) := [h(\hat{x}) + \langle A\hat{x}, y \rangle - J(y)] - [h(x) + \langle Ax, \hat{y} \rangle - J(\hat{y})], \quad \forall z = (x, y) \in Z. \quad (2.5)$$

It can be easily verified that  $\hat{z} \in Z$  is an optimal solution of problem (1.1)-(1.2) if and only if  $Q_0(\hat{z}, z) \leq 0$  for any  $z \in Z$ . A natural way to assess the solution quality of  $\hat{z}$  is to compute the gap

$$g_0(\hat{z}) = \max_{z \in Z} Q_0(\hat{z}, z), \quad (2.6)$$

under the assumption that  $g_0$  is well-defined, e.g., when  $Z$  is bounded [6, 9]. Since  $\hat{z}$  is a random variable in the randomized primal-dual algorithm to be studied in this paper, one would expect to use  $\mathbf{E}[g_0(\hat{z})]$  to measure the quality of  $\hat{z}$ . However, except for a few specific cases, we cannot provide an error bound on  $\mathbf{E}[g_0(\hat{z})]$  in general. Instead, we will introduce a slightly relaxed termination criterion defined as follows. For any given  $\delta \in \mathbb{R}$ , let us denote

$$Q_\delta(\hat{z}, z) := [h(\hat{x}) + \langle A\hat{x}, y \rangle - J(y)] - [h(x) + \langle Ax, \hat{y} \rangle - J(\hat{y})] + \delta, \quad \forall z = (x, y) \in Z \quad (2.7)$$

and

$$g_\delta(\hat{z}) := \max_{z \in Z} Q_\delta(\hat{z}, z). \quad (2.8)$$

We will show the convergence of the randomized primal-dual algorithm in terms of the expected primal-dual gap  $\mathbf{E}[g_\delta(\hat{z})]$  for some  $\delta \in \mathbb{R}$  satisfying  $\mathbf{E}[\delta] = 0$ . Clearly,  $g_0$  in (2.6) is a specialized version of  $g_\delta$  with  $\delta = 0$ .

One potential problem associated with the aforementioned primal-dual gap  $g_\delta$  is that it is not well-defined if  $Z$  is unbounded. In the latter case, Monteiro and Svaiter [32] suggested a perturbation-based termination criterion for solving problem (1.1)-(1.2) inspired by the enlargement of a maximal monotone operator that was first studied in [5]. One advantage of using this criterion is that its definition does not depend on the boundedness of the domain of the operator. More specifically, as shown in [32], there always exists a perturbation vector  $v$  such that

$$\tilde{g}_0(\hat{z}, v) := \max_{z \in Z} Q_0(\hat{z}, z) - \langle v, \hat{z} - z \rangle$$

is well-defined, although the value of  $g_0(\hat{z})$  in (2.6) may be unbounded if  $Z$  is unbounded. Accordingly, for the case when  $\hat{z}$  is a random variable, we define

$$\tilde{g}_\delta(\hat{z}, v) := \max_{z \in Z} Q_\delta(\hat{z}, z) - \langle v, \hat{z} - z \rangle \quad (2.9)$$

and establish the convergence of the randomized primal-dual algorithm in terms of  $\mathbf{E}[\tilde{g}_\delta(\hat{z}, v)]$  for some  $\delta \in \mathbb{R}$  satisfying  $\mathbf{E}[\delta] = 0$ .

### 3 The algorithm and main results

This section consists of three subsections. We first present a generic randomized primal-dual (RPD) method in subsection 3.1, and discuss its convergence properties for solving different classes of saddle point problems in the two subsequent subsections. More specifically, we focus on the analysis of the RPD method for solving general saddle point problems, where both  $h$  and  $J$  are general convex functions without assuming strong convexity, over bounded feasible sets in subsection 3.2. We then show in subsection 3.3 that much stronger convergence properties can be obtained for solving smooth saddle point problems, for which  $J$  is strongly convex. It is worth noting that the same algorithmic framework presented in subsection 3.1 is applicable to all these different cases mentioned above, as well as the unbounded case to be discussed in Section 4.

#### 3.1 The RPD algorithm

We will first introduce a few notations in order to simplify the description and analysis of the RPD algorithm. Let  $I_m$  and  $I_{m_i}$ ,  $i = 1, 2, \dots, p$ , respectively, denote the identity matrices in  $\mathbb{R}^{m \times m}$  and  $\mathbb{R}^{m_i \times m_i}$ ,  $i = 1, 2, \dots, p$ . Observe that  $I_{m_i}$ ,  $i = 1, 2, \dots, p$ , can be viewed as the  $i$ -th diagonal block of  $I_m$ . Also let us define  $U_i \in \mathbb{R}^{m \times m}$ ,  $i = 1, 2, \dots, p$ , as the diagonal matrix whose  $i$ -th diagonal block is  $I_{m_i}$  and all other blocks are given by 0. Also let  $\bar{U}_i \in \mathbb{R}^{m \times m}$  be the complement of  $U_i$  such that

$$U_i + \bar{U}_i = I_m.$$

With the help of the above notations, we are now ready to describe our algorithmic framework as follows.

---

**Algorithm 1** The randomized primal-dual (RPD) method for saddle point optimization

---

Let  $z^1 = (x^1, y^1) \in X \times Y$ , and nonnegative stepsizes  $\{\tau_t\}$ ,  $\{\eta_t\}$ , parameters  $\{q_t\}$ , and weights  $\{\gamma_t\}$  be given. Set  $\bar{x}^1 = x^1$ .

**for**  $t = 1, \dots, N$  **do**

1. Generate a random variable  $i_t$  uniformly distributed over  $\{1, 2, \dots, p\}$ .
2. Update  $y^{t+1}$  and  $x^{t+1}$  by

$$y_i^{t+1} = \begin{cases} \operatorname{argmin}_{y_i \in Y_i} \langle -U_i A \bar{x}^t, y \rangle + J_i(y_i) + \frac{\tau_t}{2} \|y_i - y_i^t\|_2^2, & i = i_t, \\ y_i^t, & i \neq i_t. \end{cases} \quad (3.10)$$

$$x^{t+1} = \operatorname{argmin}_{x \in X} h(x) + \langle x, A^T y^{t+1} \rangle + \frac{\eta_t}{2} \|x - x^t\|_2^2. \quad (3.11)$$

$$\bar{x}^{t+1} = q_t(x^{t+1} - x^t) + x^{t+1}. \quad (3.12)$$

**end for**

**Output:** Set

$$\hat{z}^N = \left( \sum_{t=1}^{N-1} \gamma_t \right)^{-1} \sum_{t=1}^{N-1} \gamma_t z^{t+1}. \quad (3.13)$$


---

The above RPD algorithm originated from the primal-dual method in [6]. The major differences between these two algorithms are summarized as follows. Firstly, instead of updating the whole dual variable  $y_i^t$ ,  $i = 1, \dots, p$ , as in the original primal-dual algorithm, the RPD algorithm updates in Step (3.10) the  $i_t$ -th component of  $y^t$  only. Secondly, rather than using constant stepsizes for  $\tau_t$ ,  $\eta_t$ , and  $q_t$ , variable stepsizes are used in the RPD method. Thirdly, the output solution  $\hat{z}^N$  is defined as a weighted average rather than a simple average of  $z^t$ ,  $t = 2, \dots, N + 1$ . The latter two enhancements are introduced so that the primal-dual algorithm can achieve the optimal rate of convergence for solving smooth saddle point problems, which is new even for the deterministic case where the number of blocks  $p = 1$ .

It is also known that the primal-dual algorithm is related to the Douglas-Rachford splitting method [11, 29] and a pre-conditioned version of the alternating direction method of multipliers [14, 17] (see, e.g., [3, 6, 13, 21, 35] for detailed reviews on the relationship between the primal-dual methods and other algorithms, as well as recent theoretical developments). However, to the best of our knowledge, there does not exist randomized version of these algorithms which only need to solve one dual subproblem at each iteration before in the literature (see Section 4 for more discussions).

It should be noted that Algorithm 1 is conceptual only since we have not yet specified a few algorithmic parameters including  $\{\tau_t\}$ ,  $\{\eta_t\}$ ,  $\{q_t\}$ , and  $\{\gamma_t\}$ . We will come back to this issue after establishing some convergence properties of the generic RPD method for solving different classes of saddle-point problems.

### 3.2 General bilinear saddle point problems over bounded feasible sets

Throughout this subsection we assume that both  $h$  and  $J$  are general convex function (without assuming strong convexity) so that problems (3.10) and (3.11) are relatively easy to solve. Also we assume that both  $X$  and  $Y$  are bounded, i.e.,  $\exists \Omega_X > 0$  and  $\Omega_Y > 0$  such that

$$\max_{x_1, x_2 \in X} \|x_1 - x_2\|_2^2 \leq \Omega_X^2 \text{ and } \max_{y_1, y_2 \in Y} \|y_1 - y_2\|_2^2 \leq \Omega_Y^2. \quad (3.14)$$

Before establishing the main convergence properties for the RPD method applied to general bilinear saddle point problems, we show an important recursion of this algorithm in the following result.

**Proposition 1** *Let  $z^t = (x^t, y^t)$ ,  $t = 1, 2, \dots, N$ , be generated by Algorithm 1. For any  $z \in Z$ , we have*

$$\begin{aligned} & \gamma_t Q_0(z^{t+1}, z) + \langle \gamma_t A x^{t+1} - A x^t, y^{t+1} - y \rangle + (\gamma_t - 1) [J(y) - J(y^{t+1})] - \Delta_t \\ & \leq \frac{\gamma_t \eta_t}{2} [\|x - x^t\|_2^2 - \|x^t - x^{t+1}\|_2^2 - \|x - x^{t+1}\|_2^2] \\ & \quad + \frac{\tau_t}{2} [\|y - y^t\|_2^2 - \|y - y^{t+1}\|_2^2 - \|y^t - y^{t+1}\|_2^2], \end{aligned} \quad (3.15)$$

where

$$\Delta_t := \langle q_{t-1} U_{i_t} A(x^t - x^{t-1}), y^{t+1} - y \rangle - \langle \bar{U}_{i_t} A x^t, y^t - y \rangle + \sum_{i \neq i_t} [J(y_i^t) - J_i(y_i)]. \quad (3.16)$$

*Proof.* By the optimality condition of problem (3.11), for all  $x \in X$ , we have

$$h(x^{t+1}) - h(x) + \langle x^{t+1} - x, A^T y^{t+1} \rangle + \frac{\eta_t}{2} \|x^t - x^{t+1}\|_2^2 + \frac{\eta_t}{2} \|x - x^{t+1}\|_2^2 \leq \frac{\eta_t}{2} \|x - x^t\|_2^2. \quad (3.17)$$

Observe that

$$\begin{aligned}\langle x^{t+1} - x, A^T y^{t+1} \rangle &= \langle Ax^{t+1}, y \rangle - \langle Ax, y^{t+1} \rangle - \langle Ax^{t+1}, y \rangle + \langle Ax^{t+1}, y^{t+1} \rangle \\ &= \langle Ax^{t+1}, y \rangle - \langle Ax, y^{t+1} \rangle + \langle Ax^{t+1}, y^{t+1} - y \rangle,\end{aligned}$$

which together with (3.17) and the definition of  $Q_0$  in (2.5) then imply

$$\begin{aligned}Q_0(z^{t+1}, z) + \langle Ax^{t+1}, y^{t+1} - y \rangle + J(y) - J(y^{t+1}) \\ \leq \frac{\eta_t}{2} [\|x - x^t\|_2^2 - \|x^t - x^{t+1}\|_2^2 - \|x - x^{t+1}\|_2^2].\end{aligned}\quad (3.18)$$

Now, by the optimality condition of problem (3.10), for all  $y \in Y$ , we have

$$\langle -U_{i_t} A \bar{x}^t, y^{t+1} - y \rangle + J_{i_t}(y_{i_t}^{t+1}) - J_{i_t}(y_{i_t}) + \frac{\tau_t}{2} \|y_{i_t}^t - y_{i_t}^{t+1}\|_2^2 + \frac{\tau_t}{2} \|y_{i_t} - y_{i_t}^{t+1}\|_2^2 \leq \frac{\tau_t}{2} \|y_{i_t} - y_{i_t}^t\|_2^2. \quad (3.19)$$

Using the definition of  $\bar{x}^t$  in (3.12), we also have

$$\begin{aligned}\langle -U_{i_t} A \bar{x}^t, y^{t+1} - y \rangle &= \langle -U_{i_t} A [q_{t-1}(x^t - x^{t-1}) + x^t], y^{t+1} - y \rangle \\ &= \langle -U_{i_t} A x^t, y^{t+1} - y \rangle - \langle q_{t-1} U_{i_t} A (x^t - x^{t-1}), y^{t+1} - y \rangle \\ &= \langle -(U_{i_t} + \bar{U}_{i_t}) A x^t, y^{t+1} - y \rangle - \langle q_{t-1} U_{i_t} A (x^t - x^{t-1}), y^{t+1} - y \rangle \\ &\quad + \langle \bar{U}_{i_t} A x^t, y^{t+1} - y \rangle \\ &= \langle -A x^t, y^{t+1} - y \rangle - \langle q_{t-1} U_{i_t} A (x^t - x^{t-1}), y^{t+1} - y \rangle + \langle \bar{U}_{i_t} A x^t, y^t - y \rangle,\end{aligned}\quad (3.20)$$

where the last identity follows from the fact that  $U_{i_t} + \bar{U}_{i_t} = I_n$  and that  $\langle \bar{U}_{i_t} A x^t, y^{t+1} - y \rangle = \langle \bar{U}_{i_t} A x^t, y^t - y \rangle$ . Also observe that

$$\begin{aligned}J_{i_t}(y_{i_t}^{t+1}) - J_{i_t}(y_{i_t}) &= J(y^{t+1}) - \\ \text{sum}_{i \neq i_t} J_i(y_i^t) - [J(y) - \sum_{i \neq i_t} J_i(y_i)] &= J(y^{t+1}) - J(y) - \sum_{i \neq i_t} [J(y_i^t) - J_i(y_i)], \\ \|y_{i_t}^t - y_{i_t}^{t+1}\|_2^2 &= \|y^t - y^{t+1}\|_2^2, \\ \|y_{i_t} - y_{i_t}^t\|_2^2 - \|y_{i_t} - y_{i_t}^{t+1}\|_2^2 &= \|y - y^t\|_2^2 - \|y - y^{t+1}\|_2^2.\end{aligned}\quad (3.21)$$

Using these observations in (3.19), we conclude

$$\begin{aligned}\langle -A x^t, y^{t+1} - y \rangle - \langle q_{t-1} U_{i_t} A (x^t - x^{t-1}), y^{t+1} - y \rangle + \langle \bar{U}_{i_t} A x^t, y^t - y \rangle \\ + J(y^{t+1}) - J(y) - \sum_{i \neq i_t} [J(y_i^t) - J_i(y_i)] \leq \frac{\tau_t}{2} [\|y - y^t\|_2^2 - \|y - y^{t+1}\|_2^2 - \|y^t - y^{t+1}\|_2^2].\end{aligned}$$

Multiplying both sides of (3.18) by  $\gamma_t$  and adding it up with the above inequality, we have

$$\begin{aligned}\gamma_t Q_0(z^{t+1}, z) + \langle \gamma_t A x^{t+1} - A x^t, y^{t+1} - y \rangle + (\gamma_t - 1) [J(y) - J(y^{t+1})] \\ - \langle q_{t-1} U_{i_t} A (x^t - x^{t-1}), y^{t+1} - y \rangle + \langle \bar{U}_{i_t} A x^t, y^t - y \rangle - \sum_{i \neq i_t} [J(y_i^t) - J_i(y_i)] \\ \leq \frac{\gamma_t \eta_t}{2} [\|x - x^t\|_2^2 - \|x^t - x^{t+1}\|_2^2 - \|x - x^{t+1}\|_2^2] \\ + \frac{\tau_t}{2} [\|y - y^t\|_2^2 - \|y - y^{t+1}\|_2^2 - \|y^t - y^{t+1}\|_2^2],\end{aligned}$$

which, in view of the definition of  $\Delta_t$ , clearly implies the result.  $\blacksquare$

The following lemma provides an upper bound on  $\mathbf{E}_{i_t}[\Delta_t]$ .

**Lemma 2** *Let  $\Delta_t$  be defined in (3.16). If  $i_t$  is uniformly distributed on  $\{1, 2, \dots, p\}$ , then*

$$\begin{aligned}\mathbf{E}_{i_t}[\Delta_t] \leq \left\langle \left( \frac{1}{p} q_{t-1} - \frac{p-1}{p} \right) A x^t - \frac{1}{p} q_{t-1} A x^{t-1}, y^t - y \right\rangle + \frac{p-1}{p} [J(y^t) - J(y)] \\ + \frac{q_{t-1}^2 \|A\|_2^2}{2p\tau_t} \|x^t - x^{t-1}\|_2^2 + \frac{\tau_t}{2} \mathbf{E}_{i_t} [\|y^{t+1} - y^t\|_2^2].\end{aligned}$$

*Proof.* The definition of  $\Delta_t$  in (3.16) can be rewritten as

$$\begin{aligned}\Delta_t &= \langle q_{t-1}U_{i_t}A(x^t - x^{t-1}) - \bar{U}_{i_t}Ax^t, y^t - y \rangle \\ &\quad - \langle q_{t-1}U_{i_t}A(x^t - x^{t-1}), y^{t+1} - y^t \rangle + \sum_{i \neq i_t} [J_i(y_i^t) - J_i(y_i)].\end{aligned}\quad (3.22)$$

Since  $i_t$  is uniformly distributed on  $\{1, 2, \dots, p\}$ , we have

$$\begin{aligned}\mathbf{E}_{i_t} &[\langle q_{t-1}U_{i_t}A(x^t - x^{t-1}) - \bar{U}_{i_t}Ax^t, y^t - y \rangle] \\ &= \left\langle \frac{1}{p}q_{t-1}A(x^t - x^{t-1}), y^t - y \right\rangle - \frac{p-1}{p} \langle Ax^t, y^t - y \rangle \\ &= \left\langle \left( \frac{1}{p}q_{t-1} - \frac{p-1}{p} \right) Ax^t - \frac{1}{p}q_{t-1}Ax^{t-1}, y^t - y \right\rangle\end{aligned}\quad (3.23)$$

and

$$\mathbf{E}_{i_t} \left[ \sum_{i \neq i_t} (J_i(y_i^t) - J_i(y_i)) \right] = \frac{p-1}{p} [J(y^t) - J(y)]. \quad (3.24)$$

Observe that

$$\begin{aligned}\mathbf{E}_{i_t} [\langle q_{t-1}U_{i_t}A(x^t - x^{t-1}), y^{t+1} - y^t \rangle] &\leq \mathbf{E}_{i_t} [q_{t-1} \|U_{i_t}A(x^t - x^{t-1})\|_2 \|y^{t+1} - y^t\|_2] \\ &\leq \mathbf{E}_{i_t} \left[ \frac{q_{t-1}^2}{2\tau_t} \|U_{i_t}A(x^t - x^{t-1})\|_2^2 + \frac{\tau_t}{2} \|y^{t+1} - y^t\|_2^2 \right] \\ &= \frac{q_{t-1}^2}{2p\tau_t} \|A(x^t - x^{t-1})\|_2^2 + \frac{\tau_t}{2} \mathbf{E}_{i_t} [\|y^{t+1} - y^t\|_2^2] \\ &\leq \frac{q_{t-1}^2 \|A\|_2^2}{2p\tau_t} \|x^t - x^{t-1}\|_2^2 + \frac{\tau_t}{2} \mathbf{E}_{i_t} [\|y^{t+1} - y^t\|_2^2],\end{aligned}\quad (3.25)$$

where the second inequality follows from the Cauchy-Swartz inequality and the equality follows from the fact that  $i_t$  is uniformly distributed on  $\{1, 2, \dots, p\}$ . The result immediately follows from (3.22), (3.23), (3.24), and (3.25).  $\blacksquare$

We are now ready to establish the main convergence properties of the RPD algorithm for solving saddle point problems over bounded feasible sets.

**Theorem 3** *Suppose that the initial point of Algorithm 1 is chosen such that  $x^1 = x^0$  and  $y^1 = \operatorname{argmax}_{y \in Y} \langle Ax^1, y \rangle - J(y)$ . Also assume that the parameters  $\{q_t\}$ ,  $\{\gamma_t\}$ ,  $\{\tau_t\}$ , and  $\{\eta_t\}$  satisfy*

$$q_t = p, \quad t = 1, \dots, N-1, \quad (3.26)$$

$$\gamma_t = \frac{1}{p}q_t - \frac{p-1}{p}, \quad t = 1, \dots, N-2 \text{ and } \gamma_{N-1} = 1, \quad (3.27)$$

$$\tau_{t-1} \geq \tau_t, \quad i = 1, \dots, N-1, \quad (3.28)$$

$$\gamma_{t-1}\eta_{t-1} \geq \gamma_t\eta_t, \quad i = 1, \dots, N-1, \quad (3.29)$$

$$p\gamma_t\eta_t\tau_{t+1} \geq q_t^2 \|A\|_2^2, \quad i = 1, \dots, N-2, \quad (3.30)$$

$$\gamma_{N-1}\eta_{N-1}\tau_{N-1} \geq \|A\|_2^2. \quad (3.31)$$

a) For any  $N \geq 1$ , we have

$$\mathbf{E}[Q_0(\hat{z}^N, z)] \leq \left( \sum_{t=1}^{N-1} \gamma_t \right)^{-1} \left[ \frac{\gamma_1\eta_1}{2} \Omega_X^2 + \frac{\tau_1}{2} \Omega_Y^2 \right], \quad \forall z \in Z, \quad (3.32)$$

where  $\hat{z}^N$  is defined in (3.13) and the expectation is taken w.r.t.  $[i_N] = (i_1, \dots, i_{N-1})$ .



b) For any  $N \geq 1$ , there exists a function  $\delta(y)$  such that  $\mathbb{E}[\delta(y)] = 0$  for any  $y \in Y$  and

$$\mathbf{E}[g_{\delta(y)}(\hat{z}^N)] \leq \left( \sum_{t=1}^{N-1} \gamma_t \right)^{-1} \left[ \frac{\gamma_1 \eta_1}{2} \Omega_X^2 + \frac{\tau_1}{2} \Omega_Y^2 \right], \quad \forall z \in Z. \quad (3.33)$$

*Proof.* We first show part a). It follows from Proposition 1 and Lemma 2 that

$$\begin{aligned} & \gamma_t Q_0(z^{t+1}, z) + \langle \gamma_t A x^{t+1} - A x^t, y^{t+1} - y \rangle + (\gamma_t - 1) [J(y) - J(y^{t+1})] \\ & \leq \mathbf{E}_{i_t}[\Delta_t] + \frac{\gamma_t \eta_t}{2} [\|x - x^t\|_2^2 - \|x^t - x^{t+1}\|_2^2 - \|x - x^{t+1}\|_2^2] \\ & \quad + \frac{\tau_t}{2} [\|y - y^t\|_2^2 - \|y - y^{t+1}\|_2^2 - \|y^t - y^{t+1}\|_2^2] + \Delta_t - \mathbf{E}_{i_t}[\Delta_t] \\ & \leq \left\langle \left( \frac{1}{p} q_{t-1} - \frac{p-1}{p} \right) A x^t - \frac{1}{p} q_{t-1} A x^{t-1}, y^t - y \right\rangle + \frac{p-1}{p} [J(y^t) - J(y)] \\ & \quad + \frac{q_{t-1}^2 \|A\|_2^2}{2p\tau_t} \|x^t - x^{t-1}\|_2^2 + \frac{\tau_t}{2} \mathbf{E}_{i_t} [\|y^{t+1} - y^t\|_2^2] + \frac{\gamma_t \eta_t}{2} [\|x - x^t\|_2^2 - \|x^t - x^{t+1}\|_2^2 - \|x - x^{t+1}\|_2^2] \\ & \quad + \frac{\tau_t}{2} [\|y - y^t\|_2^2 - \|y - y^{t+1}\|_2^2 - \|y^t - y^{t+1}\|_2^2] + \Delta_t - \mathbf{E}_{i_t}[\Delta_t]. \end{aligned} \quad (3.34)$$

Denoting

$$\Delta'_t := \Delta_t - \frac{\tau_t}{2} \|y^t - y^{t+1}\|_2^2,$$

we can rewrite (3.34) as

$$\begin{aligned} & \gamma_t Q_0(z^{t+1}, z) + \langle \gamma_t A x^{t+1} - A x^t, y^{t+1} - y \rangle + (1 - \gamma_t) [J(y^{t+1}) - J(y)] \\ & \leq \left\langle \left( \frac{1}{p} q_{t-1} - \frac{p-1}{p} \right) A x^t - \frac{1}{p} q_{t-1} A x^{t-1}, y^t - y \right\rangle + \frac{p-1}{p} [J(y^t) - J(y)] - \frac{\gamma_t \eta_t}{2} \|x^t - x^{t+1}\|_2^2 \\ & \quad + \frac{q_{t-1}^2 \|A\|_2^2}{2p\tau_t} \|x^t - x^{t-1}\|_2^2 + \frac{\gamma_t \eta_t}{2} [\|x - x^t\|_2^2 - \|x - x^{t+1}\|_2^2] + \frac{\tau_t}{2} [\|y - y^t\|_2^2 - \|y - y^{t+1}\|_2^2] \\ & \quad + \Delta'_t - \mathbf{E}_{i_t}[\Delta'_t]. \end{aligned} \quad (3.35)$$

Taking summation from  $t = 1$  to  $N - 1$  on both sides of the above inequality, using the assumptions in (3.26) and (3.27), and denoting  $z^{[N]} := \{(x^t, y^t)\}_{t=1}^N$  and

$$\mathcal{B}_N(z, z^{[N]}) := \sum_{t=1}^{N-1} \left[ \frac{\gamma_t \eta_t}{2} \|x - x^t\|_2^2 - \frac{\gamma_t \eta_t}{2} \|x - x^{t+1}\|_2^2 \right] + \sum_{t=1}^{N-1} \left[ \frac{\tau_t}{2} \|y - y^t\|_2^2 - \frac{\tau_t}{2} \|y - y^{t+1}\|_2^2 \right], \quad (3.36)$$

we then conclude that

$$\begin{aligned} & \sum_{t=1}^{N-1} \gamma_t Q_0(z^{t+1}, z) \\ & \leq \mathcal{B}_N(z, z^{[N]}) - \langle A x^N - A x^{N-1}, y^N - y \rangle + \left\langle \frac{1}{p} A x^1 - A x^0, y^1 - y \right\rangle + \frac{p-1}{p} [J(y^1) - J(y)] \\ & \quad + \frac{p \|A\|_2^2}{2\tau_1} \|x^1 - x^0\|_2^2 - \frac{\gamma_{N-1} \eta_{N-1}}{2} \|x^N - x^{N-1}\|_2^2 - \sum_{t=1}^{N-2} \left( \frac{\gamma_t \eta_t}{2} - \frac{q_t^2 \|A\|_2^2}{2p\tau_{t+1}} \right) \|x^{t+1} - x^t\|_2^2 \\ & \quad + \sum_{t=1}^{N-1} (\Delta'_t - \mathbf{E}_{i_t}[\Delta'_t]) \\ & \leq \mathcal{B}_N(z, z^{[N]}) - \frac{\gamma_{N-1} \eta_{N-1}}{2} \|x^N - x^{N-1}\|_2^2 - \langle A x^N - A x^{N-1}, y^N - y \rangle + \sum_{t=1}^{N-1} (\Delta'_t - \mathbf{E}_{i_t}[\Delta'_t]), \end{aligned} \quad (3.37)$$

where the second inequality follows from (3.30), and the facts that  $x^1 = x^0$  and  $y^1 = \operatorname{argmax}_{y \in Y} \langle A x^1, y \rangle - J(y)$ . Using the above conclusion, the definition of  $\hat{x}^N$  in (3.13), and the convexity of  $Q_0(\hat{z}, z)$  w.r.t.  $\hat{z}$ , we obtain

$$\begin{aligned} \left( \sum_{t=1}^{N-1} \gamma_t \right) Q_0(\hat{z}^N, z) & \leq \sum_{t=1}^{N-1} \gamma_t Q_0(z^{t+1}, z) \\ & \leq \mathcal{B}_N(z, z^{[N]}) - \frac{\gamma_{N-1} \eta_{N-1}}{2} \|x^N - x^{N-1}\|_2^2 - \langle A x^N - A x^{N-1}, y^N - y \rangle \\ & \quad + \sum_{t=1}^{N-1} (\Delta'_t - \mathbf{E}_{i_t}[\Delta'_t]), \end{aligned}$$

which, in view of the fact that

$$-\langle Ax^N - Ax^{N-1}, y^N - y \rangle \leq \frac{\|A\|_2^2}{2\tau_{N-1}} \|x^N - x^{N-1}\|_2^2 + \frac{\tau_{N-1}}{2} \|y^N - y\|_2^2, \quad (3.38)$$

then implies that

$$\begin{aligned} \left( \sum_{t=1}^{N-1} \gamma_t \right) Q_0(\hat{z}^N, z) &\leq \mathcal{B}_N(z, z^{[N]}) + \frac{\tau_{N-1}}{2} \|y^N - y\|_2^2 - \left( \frac{\gamma_{N-1}\eta_{N-1}}{2} - \frac{\|A\|_2^2}{2\tau_{N-1}} \right) \|x^N - x^{N-1}\|_2^2 \\ &\quad + \sum_{t=1}^{N-1} (\Delta'_t - \mathbf{E}_{i_t}[\Delta'_t]). \end{aligned}$$

Now it follows from (3.28), (3.29), and (3.36) that

$$\begin{aligned} &\mathcal{B}_N(z, z^{[N]}) + \frac{\tau_{N-1}}{2} \|y^N - y\|_2^2 \\ &= \frac{\gamma_1\eta_1}{2} \|x - x^1\|_2^2 - \sum_{t=1}^{N-2} \left( \frac{\gamma_t\eta_t}{2} - \frac{\gamma_{t+1}\eta_{t+1}}{2} \right) \|x - x^{t+1}\|_2^2 - \frac{\gamma_{N-1}\eta_{N-1}}{2} \|x - x^N\|_2^2 \\ &\quad + \frac{\tau_1}{2} \|y - y^1\|_2^2 - \sum_{t=1}^{N-2} \left( \frac{\tau_t}{2} - \frac{\tau_{t+1}}{2} \right) \|y - y^{t+1}\|_2^2 \\ &\leq \frac{\gamma_1\eta_1}{2} \|x - x^1\|_2^2 - \frac{\gamma_{N-1}\eta_{N-1}}{2} \|x - x^N\|_2^2 + \frac{\tau_1}{2} \|y - y^1\|_2^2 \\ &\leq \frac{\gamma_1\eta_1}{2} \Omega_X^2 + \frac{\tau_1}{2} \Omega_Y^2. \end{aligned}$$

Combining the above two relations, and noting that  $\frac{\gamma_{N-1}\eta_{N-1}}{2} \geq \frac{\|A\|_2^2}{2\tau_{N-1}}$  by (3.31), we obtain

$$\left( \sum_{t=1}^{N-1} \gamma_t \right) Q_0(\hat{z}^N, z) \leq \frac{\gamma_1\eta_1}{2} \Omega_X^2 + \frac{\tau_1}{2} \Omega_Y^2 + \sum_{t=1}^{N-1} (\Delta'_t - \mathbf{E}_{i_t}[\Delta'_t]). \quad (3.39)$$

Taking expectation w.r.t  $i_t, t = 1, 2, \dots, N-1$ , noting that  $\mathbf{E}_{i_t}[\Delta'_t - \mathbf{E}_{i_t}[\Delta'_t]] = 0$  and  $\frac{p-1}{p} \leq 1$ , we obtain

$$\mathbf{E}_{[i_N]}[Q_0(\hat{z}^N, z)] \leq \left( \sum_{t=1}^{N-1} \gamma_t \right)^{-1} \left[ \frac{\gamma_{N-1}\eta_{N-1}}{2} \Omega_X^2 + \frac{\tau_{N-1}}{2} \Omega_Y^2 \right].$$

The proof of part b) is similar to that of part a). The main idea is to break down the perturbation term  $\Delta'_t$  into two parts, one independent on  $y$  and the other depending on  $y$ . More specifically, let us denote

$$\Delta'_{t1} = \langle q_{t-1}U_{i_t}A(x^t - x^{t-1}), y^t \rangle - \langle \bar{U}_{i_t}Ax^t, y^t \rangle + \sum_{i \neq i_t} J_i(y_i^t) - \frac{\tau_t}{2} \|y^t - y^{t+1}\|_2^2, \quad (3.40)$$

$$\Delta'_{t2} = \langle q_{t-1}U_{i_t}A(x^t - x^{t-1}), y \rangle - \langle \bar{U}_{i_t}Ax^t, y \rangle + \sum_{i \neq i_t} J_i(y_i). \quad (3.41)$$

Clearly, we have

$$\Delta'_t = \Delta'_{t1} + \Delta'_{t2}. \quad (3.42)$$

Using exactly the same analysis as in part a) except putting the perturbation term  $\Delta'_{t2}$  to the left hand side of (3.39), we have

$$\left( \sum_{t=1}^{N-1} \gamma_t \right) Q_0(\hat{z}^N, z) + \sum_{t=1}^{N-1} (\Delta'_{t2} - \mathbf{E}_{i_t}[\Delta'_{t2}]) \leq \frac{\gamma_1\eta_1}{2} \Omega_X^2 + \frac{\tau_1}{2} \Omega_Y^2 + \sum_{t=1}^{N-1} (\Delta'_{t1} - \mathbf{E}_{i_t}[\Delta'_{t1}]).$$

Denoting

$$\delta(y) = \left( \sum_{t=1}^{N-1} \gamma_t \right)^{-1} \sum_{t=1}^{N-1} (\Delta'_{t2} - \mathbf{E}_{i_t}[\Delta'_{t2}]),$$

we then conclude from the above inequality that

$$\left( \sum_{t=1}^{N-1} \gamma_t \right) [Q_0(\hat{z}^N, z) + \delta(y)] \leq \frac{\gamma_1\eta_1}{2} \Omega_X^2 + \frac{\tau_1}{2} \Omega_Y^2 + \sum_{t=1}^{N-1} (\Delta'_{t1} - \mathbf{E}_{i_t}[\Delta'_{t1}]).$$

The result in (3.33) then immediately follows by maximizing both sides of the above inequality w.r.t  $z = (x, y)$ , and taking expectation w.r.t  $i_t, t = 1, 2, \dots, N - 1$ , and using the definition of  $g_\delta$  in (2.8). ■

While there are many options to specify the parameters  $\eta_t, \tau_t$ , and  $\gamma_t$  of the RPD method such that the assumptions in (3.26)-(3.31) are satisfied, below we provide a specific parameter setting which leads to an optimal rate of convergence for the RPD algorithm in terms of its dependence on  $N$ .

**Corollary 4** *Suppose that the initial point of Algorithm 1 is set to  $x^1 = x^0$  and  $y^1 = \operatorname{argmax}_{y \in Y} \langle Ax^1, y \rangle - J(y)$ . Also assume that  $q_t$  is set to (3.26), and  $\{\gamma_t\}, \{\tau_t\}$ , and  $\{\eta_t\}$  are set to*

$$\gamma_t = \frac{1}{p}, t = 1, 2, \dots, N - 2, \text{ and } \gamma_{N-1} = 1, \quad (3.43)$$

$$\tau_t = \frac{\sqrt{p}\|A\|\Omega_X}{\Omega_Y}, \quad (3.44)$$

$$\eta_t = \frac{p^{\frac{3}{2}}\|A\|\Omega_Y}{\Omega_X}, t = 1, 2, \dots, N - 2, \text{ and } \eta_{N-1} = \frac{\sqrt{p}\|A\|\Omega_Y}{\Omega_X}. \quad (3.45)$$

Then for any  $N \geq 1$ , we have

$$\mathbf{E}[Q_0(\hat{z}^N, z)] \leq \frac{p^{3/2}\|A\|_2\Omega_X\Omega_Y}{N+p-2}, \forall z \in Z, \quad (3.46)$$

Moreover, there exists a function  $\delta(y)$  such that  $\mathbb{E}[\delta(y)] = 0$  for any  $y \in Y$  and

$$\mathbf{E}[g_\delta(y)(\hat{z}^N)] \leq \frac{p^{3/2}\|A\|_2\Omega_X\Omega_Y}{N+p-2}. \quad (3.47)$$

*Proof.* It is easy to verify that  $\gamma_t, \tau_t$ , and  $\eta_t$  defined in (3.43)-(3.45) satisfy (3.27)-(3.31). Moreover, it follows from (3.43)-(3.45) that

$$\sum_{t=1}^{N-1} \gamma_t = \frac{N+p-2}{p}, \quad \frac{\gamma_1\eta_1}{2}\Omega_X^2 = \frac{\sqrt{p}\|A\|_2\Omega_X\Omega_Y}{2} \text{ and } \frac{\tau_1}{2}\Omega_Y^2 = \frac{\sqrt{p}\|A\|_2\Omega_X\Omega_Y}{2}.$$

The results then follow by plugging these identities into (3.32) and (3.33). ■

We now make some remarks about the convergence results obtained in Theorem 3 and Corollary 4. Observe that, in the view of (3.46), the total number of iterations required by the RPD algorithm to find an  $\epsilon$ -solution of problem (1.1), i.e., a point  $\hat{z} \in Z$  such that  $\mathbf{E}[Q_0(\hat{z}, z)] \leq \epsilon$  for any  $z \in Z$ , can be bounded by  $\mathcal{O}(p^{3/2}\|A\|_2\Omega_X\Omega_Y/\epsilon)$ . This bound is not improvable in terms of its dependence on  $\epsilon$  for a given  $p$  (see discussions in [9]). It should be noted, however, that the number of dual subproblems to be solved in the RPD algorithm is larger than the one required by the deterministic primal-dual method, i.e.,  $\mathcal{O}(p\|A\|_2\Omega_X\Omega_Y/\epsilon)$ , by a factor of  $\sqrt{p}$ . On the other hand, in comparison with stochastic algorithms such as the stochastic mirror descent (SMD) method (see [37, 24, 10]), Algorithm 1 exhibits a significantly better dependence on  $\epsilon$ , as the latter algorithm would require  $\mathcal{O}(1/\epsilon^2)$  iterations to find an  $\epsilon$ -solution of problem (1.1)-(1.2).

### 3.3 Smooth bilinear saddle point problems over bounded feasible sets

In this section, we assume that  $J_i(y_i), i = 1, 2, \dots, p$ , in (1.1)-(1.2) are strongly convex functions. Moreover, without loss of generality we assume that their strong convexity modulus is given by 1.

Under these assumptions, the objective function of (1.1) is a smooth convex function, which explains why these problems are called smooth bilinear saddle point problems. Our goal is to show that the RPD algorithm, when equipped with properly specified algorithmic parameters, exhibits an optimal  $\mathcal{O}(1/N^2)$  rate of convergence for solving this class of saddle point problems.

Similar to Proposition 1, we first establish an important recursion for the RPD algorithm applied to smooth bilinear saddle point problems. Note that this result involves an extra parameter  $\theta_t$  in comparison with Proposition 1.

**Proposition 5** *Let  $z^t = (x^t, y^t)$ ,  $t = 1, \dots, N$  be generated by the RPD algorithm. For any  $z \in Z$ , we have*

$$\begin{aligned} & \gamma_t Q_0(z^{t+1}, z) + \langle \gamma_t A x^{t+1} - \theta_t A x^t, y^{t+1} - y \rangle + (\theta_t - \gamma_t) [J(y^{t+1}) - J(y)] - \tilde{\Delta}_t \\ & \leq \frac{\gamma_t \eta_t}{2} [\|x - x^t\|_2^2 - \|x^t - x^{t+1}\|_2^2 - \|x - x^{t+1}\|_2^2] + \frac{\theta_t \tau_t}{2} \|y - y^t\|_2^2 \\ & \quad - \theta_t \frac{1+\tau_t}{2} \|y - y^{t+1}\|_2^2 - \frac{\theta_t \tau_t}{2} \|y^t - y^{t+1}\|_2^2 \end{aligned} \quad (3.48)$$

for any  $\theta_t \geq 0$ , where

$$\begin{aligned} \tilde{\Delta}_t & := \langle q_{t-1} \theta_t U_{i_t} A(x^t - x^{t-1}), y^{t+1} - y \rangle - \langle \theta_t \bar{U}_{i_t} A x^t, y^t - y \rangle \\ & \quad + \sum_{i \neq i_t} \theta_t [J_j(y_i^t) - J_i(y_i) + \frac{1}{2} \|y - y_i^{t+1}\|_2^2]. \end{aligned} \quad (3.49)$$

*Proof.* It follows from the optimality condition of problem (3.10) (e.g., Lemma 6 of [25] and Lemma 2 of [16]) and the strong convexity of  $J_{i_t}(y_{i_t})$  (modulus 1) that for all  $y \in Y$ ,

$$\begin{aligned} & \langle -U_{i_t} A \bar{x}^t, y^{t+1} - y \rangle + J_{i_t}(y_{i_t}^{t+1}) - J_{i_t}(y_{i_t}) \\ & \leq \frac{\tau_t}{2} \|y_{i_t} - y_{i_t}^t\|_2^2 - \frac{\tau_t}{2} \|y_{i_t}^t - y_{i_t}^{t+1}\|_2^2 - \frac{1+\tau_t}{2} \|y_{i_t} - y_{i_t}^{t+1}\|_2^2. \end{aligned} \quad (3.50)$$

This relation, in view of the observations in (3.20) and (3.21), then implies that

$$\begin{aligned} & \langle -A x^t, y^{t+1} - y \rangle - \langle q_{t-1} U_{i_t} A(x^t - x^{t-1}), y^{t+1} - y \rangle + \langle \bar{U}_{i_t} A x^t, y^t - y \rangle + J(y^{t+1}) - J(y) \\ & \leq \sum_{i \neq i_t} [J_i(y_i^t) - J_i(y_i)] + \frac{\tau_t}{2} [\|y - y^t\|_2^2 - \|y - y^{t+1}\|_2^2 - \|y^t - y^{t+1}\|_2^2] - \frac{1}{2} \|y_{i_t} - y_{i_t}^{t+1}\|_2^2 \\ & = \sum_{i \neq i_t} [J_i(y_i^t) - J_i(y_i)] + \frac{\tau_t}{2} \|y - y^t\|_2^2 - \frac{1+\tau_t}{2} \|y - y^{t+1}\|_2^2 - \frac{\tau_t}{2} \|y^t - y^{t+1}\|_2^2 + \frac{1}{2} \sum_{i \neq i_t} \|y_i - y_i^{t+1}\|_2^2, \end{aligned} \quad (3.51)$$

Multiplying both sides of the above inequality by  $\theta_t$  and both sides of (3.18) by  $\gamma_t$ , and then adding them up, we obtain

$$\begin{aligned} & \gamma_t Q_0(z^{t+1}, z) + \langle \gamma_t A x^{t+1} - \theta_t A x^t, y^{t+1} - y \rangle + (\theta_t - \gamma_t) [J(y^{t+1}) - J(y)] \\ & - \langle q_{t-1} \theta_t U_{i_t} A(x^t - x^{t-1}), y^{t+1} - y \rangle + \langle \theta_t \bar{U}_{i_t} A x^t, y^t - y \rangle \\ & \leq \sum_{i \neq i_t} \theta_t [J_i(y_i^t) - J_i(y_i)] + \frac{\gamma_t \eta_t}{2} [\|x - x^t\|_2^2 - \|x^t - x^{t+1}\|_2^2 - \|x - x^{t+1}\|_2^2] + \frac{\theta_t \tau_t}{2} \|y - y^t\|_2^2 \\ & \quad - \theta_t \left( \frac{1+\tau_t}{2} \right) \|y - y^{t+1}\|_2^2 - \frac{\theta_t \tau_t}{2} \|y^t - y^{t+1}\|_2^2 + \sum_{i \neq i_t} \frac{\theta_t}{2} \|y_i - y_i^{t+1}\|_2^2, \end{aligned}$$

which, in view of the definition of  $\tilde{\Delta}_t$  in (3.49), then implies (3.48).  $\blacksquare$

The following lemma provides an upper bound on  $\mathbf{E}_{i_t}[\tilde{\Delta}_t]$ .

**Lemma 6** *Let  $\tilde{\Delta}_t$  be defined in (3.49). If  $i_t$  is uniformly distributed on  $\{1, 2, \dots, p\}$ , then*

$$\begin{aligned} \mathbf{E}_{i_t}[\tilde{\Delta}_t] & \leq \left\langle \left( \frac{1}{p} q_{t-1} \theta_t - \frac{p-1}{p} \theta_t \right) A x^t - \frac{1}{p} q_{t-1} \theta_t A x^{t-1}, y^t - y \right\rangle + \frac{p-1}{p} \theta_t [J(y^t) - J(y)] \\ & \quad + \frac{q_{t-1}^2 \theta_t \|A\|_2^2}{2p\tau_t} \|x^t - x^{t-1}\|_2^2 + \frac{\tau_t \theta_t}{2} \mathbf{E}_{i_t} [\|y^{t+1} - y^t\|_2^2] + \frac{p-1}{2p} \theta_t \|y - y^t\|_2^2. \end{aligned}$$

*Proof.* By the definition of  $\tilde{\Delta}_t$  in (3.49), we have

$$\begin{aligned} \tilde{\Delta}_t = & \langle q_{t-1}\theta_t U_{i_t} A(x^t - x^{t-1}) - \theta_t \bar{U}_{i_t} A x^t, y^t - y \rangle - \langle q_{t-1}\theta_t U_{i_t} A(x^t - x^{t-1}), y^{t+1} - y^t \rangle \\ & + \sum_{i \neq i_t} \theta_t [J_i(y_i^t) - J_i(y_i) + \frac{1}{2} \|y_i - y_i^{t+1}\|_2^2], \end{aligned}$$

The result then immediately follows from the above identity, the relations (3.23), (3.24), (3.25), and the facts that  $\theta_t \geq 0$  and

$$\mathbf{E}_{i_t} \left[ \sum_{i \neq i_t} \theta_t \|y_i - y_i^{t+1}\|_2^2 \right] = \frac{p-1}{p} \theta_t \|y - y^t\|_2^2. \quad (3.52)$$

■

We are now ready to establish the main convergence properties of the RPD algorithm applied to smooth bilinear saddle point problems.

**Theorem 7** *Suppose that the initial point of Algorithm 1 is chosen such that  $x^1 = x^0$  and  $y^1 = \operatorname{argmax}_{y \in Y} \langle Ax^1, y \rangle - J(y)$ . Also assume that the parameters  $q_t$  and the weights  $\gamma_t, \theta_t$  are set to*

$$\theta_t = \frac{1}{p} q_t \theta_{t+1}, \quad i = 1, \dots, N-1, \quad (3.53)$$

$$\gamma_t = \left( \frac{1}{p} q_t - \frac{p-1}{p} \right) \theta_{t+1}, \quad t = 1, \dots, N-2, \quad (3.54)$$

$$\gamma_{N-1} = \theta_{N-1}, \quad (3.55)$$

$$\theta_t (1 + \tau_t) \geq \theta_{t+1} \left( \frac{p-1}{p} + \tau_{t+1} \right), \quad i = 1, \dots, N-1, \quad (3.56)$$

$$\gamma_t \eta_t \geq \gamma_{t+1} \eta_{t+1} \quad i = 1, \dots, N-1, \quad (3.57)$$

$$p \gamma_{t-1} \eta_{t-1} \tau_t \geq q_{t-1}^2 \theta_t \|A\|_2^2, \quad i = 1, \dots, N-1, \quad (3.58)$$

$$\eta_{N-1} \tau_{N-1} \geq \|A\|_2^2. \quad (3.59)$$

a) For any  $N \geq 1$ , we have

$$\mathbf{E}[Q_0(\hat{z}^N, z)] \leq \left( \sum_{t=1}^{N-1} \gamma_t \right)^{-1} \left[ \frac{\gamma_1 \eta_1}{2} \Omega_X^2 + \theta_1 \left( \frac{p-1}{p} + \tau_1 \right) \Omega_Y^2 \right], \quad (3.60)$$

where the expectation is taken w.r.t. to  $i_{[N]} = (i_1, \dots, i_{N-1})$ .

b) For any  $N \geq 1$ , there exists a function  $\delta(y)$  such that  $\mathbb{E}[\delta(y)] = 0$  for any  $y \in Y$  and

$$\mathbf{E}[g_{\delta(y)}(\hat{z}^N)] \leq \left( \sum_{t=1}^{N-1} \gamma_t \right)^{-1} \left[ \frac{\gamma_1 \eta_1}{2} \Omega_X^2 + \theta_1 \left( \frac{p-1}{p} + \tau_1 \right) \Omega_Y^2 \right]. \quad (3.61)$$

*Proof.* We first show part a). It follows from Proposition 5 and Lemma 6 that

$$\begin{aligned} & \gamma_t Q_0(z^{t+1}, z) + \langle \gamma_t A x^{t+1} - \theta_t A x^t, y^{t+1} - y \rangle + (\theta_t - \gamma_t) [J(y^{t+1}) - J(y)] \\ & \leq \left\langle \left( \frac{1}{p} q_{t-1} - \frac{p-1}{p} \right) \theta_t A x^t - \frac{1}{p} q_{t-1} \theta_t A x^{t-1}, y^t - y \right\rangle + \frac{p-1}{p} \theta_t [J(y^t) - J(y)] + \tilde{\Delta}_t - \mathbf{E}_{i_t}[\tilde{\Delta}_t] \\ & \quad + \frac{q_{t-1}^2 \theta_t \|A\|_2^2}{2p\tau_t} \|x^t - x^{t-1}\|_2^2 + \frac{\theta_t \tau_t}{2} \mathbf{E}_{i_t} [\|y^{t+1} - y^t\|_2^2] + \theta_t \left( \frac{p-1}{2p} + \frac{\tau_t}{2} \right) \|y - y^t\|_2^2 \\ & \quad + \frac{\gamma_t \eta_t}{2} [\|x - x^t\|_2^2 - \|x^t - x^{t+1}\|_2^2 - \|x - x^{t+1}\|_2^2] - \theta_t \left( \frac{1+\tau_t}{2} \right) \|y - y^{t+1}\|_2^2 - \frac{\theta_t \tau_t}{2} \|y^t - y^{t+1}\|_2^2. \end{aligned}$$

Denoting

$$\tilde{\Delta}'_t = \tilde{\Delta}_t - \frac{\theta_t \tau_t}{2} \|y^t - y^{t+1}\|_2^2,$$

we can rewrite the above inequality as

$$\begin{aligned}
& \gamma_t Q_0(z^{t+1}, z) + \langle \gamma_t A x^{t+1} - \theta_t A x^t, y^{t+1} - y \rangle + (\theta_t - \gamma_t) [J(y^{t+1}) - J(y)] \\
& \leq \left\langle \left( \frac{1}{p} q_{t-1} - \frac{p-1}{p} \right) \theta_t A x^t - \frac{1}{p} q_{t-1} \theta_t A x^{t-1}, y^t - y \right\rangle + \frac{p-1}{p} \theta_t [J(y^t) - J(y)] \\
& \quad + \frac{q_{t-1}^2 \theta_t \|A\|_2^2}{2p\tau_t} \|x^t - x^{t-1}\|_2^2 + \frac{\gamma_t \eta_t}{2} [\|x - x^t\|_2^2 - \|x^t - x^{t+1}\|_2^2 - \|x - x^{t+1}\|_2^2] \\
& \quad + \theta_t \left( \frac{p-1}{2p} + \frac{\tau_t}{2} \right) \|y - y^t\|_2^2 - \theta_t \frac{1+\tau_t}{2} \|y - y^{t+1}\|_2^2 + \tilde{\Delta}'_t - \mathbf{E}_{i_t}[\tilde{\Delta}'_t].
\end{aligned} \tag{3.62}$$

Observe that by (3.53), (3.54), and (3.55), and the fact  $x^1 = x^0$ ,

$$\begin{aligned}
& \langle \gamma_t A x^{t+1} - \theta_t A x^t, y^{t+1} - y \rangle = \langle \left( \frac{1}{p} q_t - \frac{p-1}{p} \right) \theta_{t+1} A x^{t+1} - \frac{1}{p} q_t \theta_{t+1} A x^t, y^{t+1} - y \rangle \\
& \langle \left( \frac{1}{p} q_0 - \frac{p-1}{p} \right) \theta_1 A x^1 - \frac{1}{p} q_0 \theta_1 A x^0, y^1 - y \rangle = \langle \left( -\frac{p-1}{p} \right) \theta_1 A x^1, y^1 - y \rangle \\
& (\theta_t - \gamma_t) [J(y^{t+1}) - J(y)] = \frac{p-1}{p} \theta_{t+1} [J(y^{t+1}) - J(y)] \\
& (\theta_{N-1} - \gamma_{N-1}) [J(y^N) - J(y)] = 0.
\end{aligned}$$

Taking summation from  $t = 1$  to  $N - 1$  on both sides of (3.62), using the above observations, and denoting

$$\begin{aligned}
\tilde{\mathcal{B}}_N(z, z^{[N]}) & := \sum_{t=1}^{N-1} \left[ \frac{\gamma_t \eta_t}{2} \|x - x^t\|_2^2 - \frac{\gamma_t \eta_t}{2} \|x - x^{t+1}\|_2^2 \right] \\
& \quad + \sum_{t=1}^{N-1} \left[ \frac{\theta_t}{2} \left( \frac{p-1}{p} + \tau_t \right) \|y - y^t\|_2^2 - \frac{\theta_t}{2} (1 + \tau_t) \|y - y^{t+1}\|_2^2 \right].
\end{aligned} \tag{3.63}$$

we obtain

$$\begin{aligned}
\sum_{t=1}^{N-1} \gamma_t Q_0(z^{t+1}, z) & \leq \tilde{\mathcal{B}}_N(z, z^{[N]}) - \langle \gamma_{N-1} A x^N - \theta_{N-1} A x^{N-1}, y^N - y \rangle - \frac{p-1}{p} \theta_1 \langle A x^1, y^1 - y \rangle \\
& \quad + \frac{p-1}{p} \theta_1 [J(y^1) - J(y)] - \sum_{t=1}^{N-2} \left( \frac{\gamma_t \eta_t}{2} - \frac{q_t^2 \theta_{t+1} \|A\|_2^2}{2p\tau_{t+1}} \right) \|x^{t+1} - x^t\|_2^2 \\
& \quad - \frac{\gamma_{N-1} \eta_{N-1}}{2} \|x^N - x^{N-1}\|_2^2 + \sum_{t=1}^{N-1} \left( \tilde{\Delta}'_t - \mathbf{E}_{i_t}[\tilde{\Delta}'_t] \right) \\
& \leq \tilde{\mathcal{B}}_N(z, z^{[N]}) - \gamma_{N-1} \langle A x^N - A x^{N-1}, y^N - y \rangle - \frac{\gamma_{N-1} \eta_{N-1}}{2} \|x^N - x^{N-1}\|_2^2 \\
& \quad + \sum_{t=1}^{N-1} \left( \tilde{\Delta}'_t - \mathbf{E}_{i_t}[\tilde{\Delta}'_t] \right),
\end{aligned} \tag{3.64}$$

where the second inequality follows from the facts that  $\gamma_{N-1} = \theta_{N-1}$  and  $y^1 = \operatorname{argmax}_{y \in Y} \langle A x^1, y \rangle - J(y)$ , and relation (3.58). The above conclusion, in view of the definition of  $\hat{x}^N$  and the convexity of  $Q_0(\hat{z}, z)$  in terms of  $\hat{z}$ , then implies that

$$\begin{aligned}
\left( \sum_{t=1}^{N-1} \gamma_t \right) Q_0(\hat{z}^N, z) & \leq \sum_{t=1}^{N-1} \gamma_t Q_0(z^{t+1}, z) \leq \tilde{\mathcal{B}}_N(z, z^{[N]}) - \frac{\gamma_{N-1} \eta_{N-1}}{2} \|x^N - x^{N-1}\|_2^2 \\
& \quad - \gamma_{N-1} \langle A x^N - A x^{N-1}, y^N - y \rangle + \sum_{t=1}^{N-1} \left( \tilde{\Delta}'_t - \mathbf{E}_{i_t}[\tilde{\Delta}'_t] \right).
\end{aligned}$$

Now by the Cauchy-Swartz inequality and the relation that  $\gamma_{N-1} = \theta_{N-1}$  in (3.55),

$$-\gamma_{N-1} \langle A x^N - A x^{N-1}, y^N - y \rangle \leq \frac{\gamma_{N-1} \|A\|_2^2}{2\tau_{N-1}} \|x^N - x^{N-1}\|_2^2 + \frac{\theta_{N-1} \tau_{N-1}}{2} \|y^N - y\|_2^2. \tag{3.65}$$

Moreover, by (3.56) and (3.57), we have

$$\begin{aligned}
\tilde{\mathcal{B}}_N(z, z^{[N]}) & = \frac{\gamma_1 \eta_1}{2} \|x - x^1\|_2^2 - \sum_{t=1}^{N-2} \left( \frac{\gamma_t \eta_t}{2} - \frac{\gamma_{t+1} \eta_{t+1}}{2} \right) \|x - x^{t+1}\|_2^2 - \frac{\gamma_{N-1} \eta_{N-1}}{2} \|x - x^N\|_2^2 \\
& \quad + \frac{\theta_1}{2} \left( \frac{p-1}{p} + \tau_1 \right) \|y - y^1\|_2^2 - \sum_{t=1}^{N-2} \left[ \frac{\theta_t}{2} (1 + \tau_t) - \frac{\theta_{t+1}}{2} \left( \frac{p-1}{p} + \tau_{t+1} \right) \right] \|y - y^{t+1}\|_2^2 \\
& \quad - \frac{\theta_{N-1}}{2} (1 + \tau_{N-1}) \|y^N - y\|_2^2 \\
& \leq \frac{\gamma_1 \eta_1}{2} \|x - x^1\|_2^2 + \frac{\theta_1}{2} \left( \frac{p-1}{p} + \tau_1 \right) \|y - y^1\|_2^2 - \frac{\theta_{N-1} \tau_{N-1}}{2} \|y^N - y\|_2^2 \\
& \leq \frac{\gamma_1 \eta_1}{2} \Omega_X^2 + \frac{\theta_1}{2} \left( \frac{p-1}{p} + \tau_1 \right) \Omega_Y^2 - \frac{\theta_{N-1} \tau_{N-1}}{2} \|y^N - y\|_2^2.
\end{aligned}$$

Combining the above three relations, and noting that  $\eta_{N-1}\tau_{N-1} \geq \|A\|_2^2$  by (3.59), we obtain

$$\left(\sum_{t=1}^{N-1} \gamma_t\right) Q_0(\hat{z}^N, z) \leq \frac{\gamma_1 \eta_1}{2} \Omega_X^2 + \theta_1 \left(\frac{p-1}{p} + \frac{\tau_1}{2}\right) \Omega_Y^2 + \sum_{t=1}^{N-1} (\tilde{\Delta}'_t - \mathbf{E}_{i_t}[\tilde{\Delta}'_t]).$$

Taking expectation w.r.t  $i_t, t = 1, 2, \dots, N-1$ , and noting that  $\mathbf{E}_{i_t}[\tilde{\Delta}'_t - \mathbf{E}_{i_t}[\tilde{\Delta}'_t]] = 0$  and  $\frac{p-1}{p} \leq 1$ , we obtain the result in part a). The proof of part b) is similar to that in Theorem 3 and hence the details are skipped.  $\blacksquare$

While there are many options to specify the parameters  $\eta_t, \tau_t, \theta_t$  and  $\gamma_t$  of the RPD method such that the assumptions in (3.53)-(3.59) are satisfied, below we provide a specific parameter setting which leads to an optimal rate of convergence for the RPD algorithm in terms of its dependence on  $N$ .

**Corollary 8** *Suppose that the initial point of Algorithm 1 is set to  $x^1 = x^0$  and  $y^1 = \operatorname{argmax}_{y \in Y} \langle Ax^1, y \rangle - J(y)$ . Also assume that  $q_t, \gamma_t, \theta_t, \tau_t$  and  $\eta_t$  are set to*

$$q_t = p \frac{t+3p}{t+3p+1}, \quad i = 1, \dots, N-1, \quad (3.66)$$

$$\gamma_t = \frac{t+2p+1}{p}, \quad t = 1, 2, \dots, N-2 \text{ and } \gamma_{N-1} = N+3p-1, \quad (3.67)$$

$$\theta_t = t+3p, \quad t = 1, 2, \dots, N-1, \quad (3.68)$$

$$\tau_t = \frac{t+p}{2p} \quad i = 1, \dots, N-1, \quad (3.69)$$

$$\eta_t = \frac{2p^3 \|A\|_2^2}{t+2p+1}. \quad (3.70)$$

Then, for any  $N \geq 1$  we have

$$\mathbf{E}_{[i_N]}[Q_0(\hat{z}^N, z)] \leq \frac{2}{N(N+p)} [p^3 \|A\|_2^2 \Omega_X^2 + 4.5p^2 \Omega_Y^2]. \quad (3.71)$$

Moreover, there exists a function  $\delta(y)$  such that  $\mathbb{E}[\delta(y)] = 0$  for any  $y \in Y$  and

$$\mathbf{E}[g_{\delta(y)}(\hat{z}^N)] \leq \frac{2}{N(N+p)} [p^3 \|A\|_2^2 \Omega_X^2 + 4.5p^2 \Omega_Y^2]. \quad (3.72)$$

*Proof.* It is easy to verify that  $q_t, \theta_t, \gamma_t, \tau_t$  and  $\eta_t$  defined in (3.66)-(3.70) satisfy (3.53)-(3.59). Moreover, it follows from (3.66)-(3.70) that

$$\sum_{t=1}^{N-1} \gamma_t = N+3p-1 + \sum_{t=1}^{N-2} \frac{t+2p+1}{p} \geq \frac{N(N+p)}{2p},$$

$$\frac{\gamma_1 \eta_1}{2} = p^2 \|A\|_2^2 \quad \text{and} \quad \theta_1 \left(\frac{p-1}{p} + \frac{\tau_1}{2}\right) = \frac{(1+3p)(3p-1)}{2p}.$$

The results then follow by plugging these identities into (3.60) and (3.61) and using the fact that  $\frac{(1+3p)(3p-1)}{2p} \leq 4.5p$ .  $\blacksquare$

We now make some remarks about the convergence results obtained in Theorem 7 and Corollary 8. Observe that, in view of (3.71), the total number of iterations required by the RPD algorithm to find an  $\epsilon$ -solution of smooth bilinear saddle point problems, i.e., a point  $\hat{z} \in Z$  such that  $\mathbf{E}[Q_0(\hat{z}, z)] \leq \epsilon$  for any  $z \in Z$ , can be bounded by

$$\max \left\{ \frac{\sqrt{2} p^{\frac{3}{2}} \|A\|_2 \Omega_X}{\sqrt{\epsilon}}, \frac{3p \Omega_Y}{\sqrt{\epsilon}} \right\}.$$

Similar to the previous results for general bilinear saddle point problems, this bound is not improvable in terms of its dependence on  $\epsilon$  for a given  $p$ .

## 4 Generalization of the randomized primal-dual method

In this section, we discuss two possible ways to generalize the RPD method. One is to extend it for solving unbounded saddle point problems and the other is to incorporate non-Euclidean distances.

### 4.1 RPD for unbounded saddle point problems

In this subsection, we assume that either the primal feasible set  $X$  or dual feasible set  $Y$  is unbounded. To assess the quality of a feasible solution  $\hat{z} \in X \times Y$ , we use the perturbation-based criterion defined in (2.9). Throughout this subsection we assume that both  $h$  and  $J$  are general convex function (without assuming strong convexity) so that problems (3.10) and (3.11) are relatively easy to solve. Our goal is to show that the RPD algorithm, when equipped with properly specified algorithmic parameters, exhibits an  $\mathcal{O}(1/N)$  rate of convergence for solving this class of unbounded saddle point problems.

Before establishing the main convergence properties for the RPD algorithm applied to unbounded bilinear saddle point problems, we first show an important property of the RPD method which states that, for every  $t \leq N$ , the expected distance from  $z^t$  to a given saddle point  $z^*$  is bounded.

**Lemma 9** *Let  $z^t = (x^t, y^t), t = 1, 2, \dots, N$ , be generated by the Algorithm 1 with  $x^1 = x^0$  and  $y^1 = \operatorname{argmax}_{y \in Y} \langle Ax^1, y \rangle - J(y)$ . Also assume that  $q_t, \tau_t, \eta_t$  and  $\gamma_t$  are set to (3.26)-(3.31). If*

$$\begin{aligned} \tau_{t-1} &= \tau_t, \quad i = 1, \dots, N-1, \\ \gamma_{t-1}\eta_{t-1} &= \gamma_t\eta_t, \quad i = 1, \dots, N-1, \\ \gamma_{t-1} &= \gamma_t, \quad i = 1, \dots, N-2, \end{aligned} \tag{4.73}$$

then we have

$$\mathbf{E}_{[i_t]} [\|x^* - x^t\|_2^2] \leq 2D^2, \quad \forall t \leq N-1, \tag{4.74}$$

$$\mathbf{E}_{[i_t]} [\|y^* - y^t\|_2^2] \leq \frac{2(2-\gamma_{t-1})\eta_{t-1}}{\tau_{t-1}} D^2, \quad \forall t \leq N-1, \tag{4.75}$$

and

$$\mathbf{E}_{[i_N]} [\|x^* - x^N\|_2^2] \leq D^2, \tag{4.76}$$

$$\mathbf{E}_{[i_N]} [\|y^* - y^N\|_2^2] \leq \frac{\eta_{N-1}\gamma_{N-1}}{\tau_{N-1}} D^2, \tag{4.77}$$

where  $[i_t] = \{i_1, \dots, i_{t-1}\}$ ,

$$D := \sqrt{\|x^* - x^1\|_2^2 + \frac{\tau_1}{\eta_1\gamma_1} \|y^* - y^1\|_2^2}, \tag{4.78}$$

and  $z^* = (x^*, y^*)$  is a saddle point of problem (1.1).

*Proof.* We first prove (4.76) and (4.77). Using (3.37) (with  $z = z^*$ ), (4.73), and the fact that  $x^1 = x^0$  and  $y^1 = \operatorname{argmax}_{y \in Y} \langle Ax^1, y \rangle - J(y)$ , we obtain

$$\begin{aligned} \sum_{t=1}^{N-1} \gamma_t Q_0(z^{t+1}, z^*) &\leq \frac{\gamma_1\eta_1}{2} \|x^1 - x^*\|_2^2 - \frac{\gamma_{N-1}\eta_{N-1}}{2} \|x^N - x^*\|_2^2 + \frac{\tau_1}{2} \|y^1 - y^*\|_2^2 - \frac{\tau_{N-1}}{2} \|y^N - y^*\|_2^2 \\ &\quad + \sum_{t=1}^{N-1} (\Delta'_t - \mathbf{E}_{i_t}[\Delta'_t]). \end{aligned}$$

Taking expectation on both sides of the above inequality w.r.t  $[i_N]$ , noting that  $Q_0(z^{t+1}, z^*) \geq 0, \forall t \geq 1$  and  $\mathbf{E}_{i_t}[\Delta'_t - \mathbf{E}_{i_t}[\Delta'_t]] = 0$ , then we obtain (4.76) and (4.77).



Now let us show (4.74) and (4.75). Using similar analysis to (3.35), we can show that for any  $j \leq t$  and any  $t \leq N - 1$ ,

$$\begin{aligned} & \gamma_j Q_0(z^{j+1}, z) + \langle \gamma_j Ax^{j+1} - Ax^j, y^{j+1} - y \rangle + (1 - \gamma_j) [J(y^{j+1}) - J(y)] \\ & \leq \left\langle \left( \frac{1}{p} q_{j-1} - \frac{p-1}{p} \right) Ax^j - \frac{1}{p} q_{j-1} Ax^{j-1}, y^j - y \right\rangle + \frac{p-1}{p} [J(y^j) - J(y)] - \frac{\gamma_t \eta_t}{2} \|x^j - x^{j+1}\|_2^2 \\ & \quad + \frac{q_{j-1}^2 \|A\|_2^2}{2p\tau_j} \|x^j - x^{j-1}\|_2^2 + \frac{\gamma_j \eta_j}{2} [\|x - x^j\|_2^2 - \|x - x^{j+1}\|_2^2] + \frac{\tau_j}{2} [\|y - y^j\|_2^2 - \|y - y^{j+1}\|_2^2] + \Delta'_j - \mathbf{E}_{i_j}[\Delta'_j]. \end{aligned}$$

Taking summation on both sides of the above inequality from  $j = 1$  to  $t - 1$  and using the facts that  $x^1 = x^0, y^1 = \operatorname{argmax}_{y \in Y} \langle Ax^1, y \rangle - J(y)$  and  $\gamma_j = \frac{1}{p}, \forall j = 1, 2, \dots, t - 1$ , we have

$$\begin{aligned} & \sum_{j=1}^{t-1} \gamma_j Q_0(z^{j+1}, z) + (1 - \gamma_{t-1}) [J(y^t) - J(y)] \\ & \leq \mathcal{B}_t(z, z^{[t]}) - \frac{\gamma_{t-1} \eta_{t-1}}{2} \|x^t - x^{t-1}\|_2^2 - \langle \gamma_{t-1} Ax^t - Ax^{t-1}, y^t - y \rangle + \sum_{j=1}^{t-1} \left( \Delta'_j - \mathbf{E}_{i_j}[\Delta'_j] \right). \end{aligned} \quad (4.79)$$

Observe that

$$\begin{aligned} -\langle \gamma_{t-1} Ax^t - Ax^{t-1}, y^t - y \rangle & = -\langle Ax^t - Ax^{t-1}, y^t - y \rangle + \langle (1 - \gamma_{t-1}) Ax^t, y^t - y \rangle \\ & = (1 - \gamma_{t-1}) [\langle Ax, y^t \rangle - \langle Ax^t, y \rangle + \langle Ax^t, y^t \rangle - \langle Ax, y^t \rangle] \\ & \quad - \langle Ax^t - Ax^{t-1}, y^t - y \rangle, \end{aligned} \quad (4.80)$$

which, in view of the fact that by the optimality condition of problem (3.11),

$$\langle Ax^t - Ax, y^t \rangle \leq \frac{\eta_{t-1}}{2} [\|x - x^{t-1}\|_2^2 - \|x - x^t\|_2^2 - \|x^t - x^{t-1}\|_2^2] - [h(x^t) - h(x)],$$

then implies that

$$\begin{aligned} -\langle \gamma_{t-1} Ax^t - Ax^{t-1}, y^t - y \rangle & \leq -\langle Ax^t - Ax^{t-1}, y^t - y \rangle + (1 - \gamma_{t-1}) [\langle Ax, y^t \rangle - \langle Ax^t, y \rangle - h(x^t) + h(x)] \\ & \quad + \frac{\eta_{t-1}(1 - \gamma_{t-1})}{2} [\|x - x^{t-1}\|_2^2 - \|x - x^t\|_2^2 - \|x^t - x^{t-1}\|_2^2]. \end{aligned}$$

By the above inequality and (4.79), we have

$$\begin{aligned} & \sum_{j=1}^{t-1} \gamma_j Q_0(z^{j+1}, z) + (1 - \gamma_{t-1}) [\langle Ax^t, y \rangle - \langle Ax, y^t \rangle + h(x^t) - h(x) + J(y^t) - J(y)] \\ & \leq \mathcal{B}_t(z, z^{[t]}) - \frac{\gamma_{t-1} \eta_{t-1}}{2} \|x^t - x^{t-1}\|_2^2 - \langle Ax^t - Ax^{t-1}, y^t - y \rangle \\ & \quad + \frac{\eta_{t-1}(1 - \gamma_{t-1})}{2} [\|x - x^{t-1}\|_2^2 - \|x - x^t\|_2^2 - \|x^t - x^{t-1}\|_2^2] + \sum_{j=1}^{t-1} \left( \Delta'_j - \mathbf{E}_{i_j}[\Delta'_j] \right). \end{aligned}$$

Using the previous relation, the fact that

$$\begin{aligned} -\langle Ax^t - Ax^{t-1}, y^t - y \rangle & \leq \|A\| \|x^t - x^{t-1}\|_2 \|y^t - y\|_2 \\ & \leq \frac{\|A\|_2^2}{\tau_{t-1}} \|x^t - x^{t-1}\|_2^2 + \frac{\tau_{t-1}}{4} \|y^t - y\|_2^2. \end{aligned}$$

and the definition of  $Q_0$  (with  $z = z^*$ ), we conclude that

$$\begin{aligned} & \sum_{j=1}^{t-1} \tilde{\gamma}_j Q_0(z^{j+1}, z^*) \\ & \leq \frac{\gamma_1 \eta_1}{2} \|x^* - x^1\|_2^2 - \left( \frac{\gamma_{t-1} \eta_{t-1}}{2} + \frac{(1 - \gamma_{t-1}) \eta_{t-1}}{2} \right) \|x^* - x^t\|_2^2 \\ & \quad - \left( \frac{\gamma_{t-1} \eta_{t-1}}{2} + \frac{(1 - \gamma_{t-1}) \eta_{t-1}}{2} - \frac{\|A\|_2^2}{\tau_{t-1}} \right) \|x^t - x^{t-1}\|_2^2 + \frac{\tau_1}{2} \|y^* - y^1\|_2^2 - \left( \frac{\tau_{t-1}}{2} - \frac{\tau_{t-1}}{4} \right) \|y^* - y^t\|_2^2 \\ & \quad + \frac{(1 - \gamma_{t-1}) \eta_{t-1}}{2} \|\hat{x} - x^{t-1}\|_2^2 + \sum_{j=1}^{t-1} \left( \Delta'_j - E_{i_j}[\Delta'_j] \right), \end{aligned}$$

where  $\tilde{\gamma}_j = \gamma_j, j = 1, \dots, t-2$  and  $\tilde{\gamma}_{t-1} = 1$ .

Now by the definition of gap function, we know that  $Q_0(z^{j+1}, z^*) \geq 0 \forall j \geq 1$ . Taking expectation on both sides of the above inequality w.r.t  $[i_t]$ , noting that  $\mathbf{E}_{i_j}[\Delta'_j - \mathbf{E}_{i_j}[\Delta'_j]] = 0$ ,  $\gamma_1\eta_1 = \gamma_{t-1}\eta_{t-1}$ ,  $\tau_1 = \tau_{t-1}$  and  $\frac{\gamma_{t-1}\eta_{t-1}}{2} + \frac{(1-\gamma_{t-1})\eta_{t-1}}{2} - \frac{\|A\|_2^2}{\tau_{t-1}} \geq 0$ , we have

$$\begin{aligned} & \left( \frac{\gamma_{t-1}\eta_{t-1}}{2} + \frac{(1-\gamma_{t-1})\eta_{t-1}}{2} \right) \mathbf{E}_{[i_t]} [\|x^* - x^t\|_2^2] + \frac{\tau_{t-1}}{4} \mathbf{E}_{[i_t]} [\|y^* - y^t\|_2^2] \\ & \leq \frac{\gamma_{t-1}\eta_{t-1}}{2} \|x^* - x^1\|_2^2 + \frac{\tau_{t-1}}{2} \|y^* - y^1\|_2^2 + \frac{(1-\gamma_{t-1})\eta_{t-1}}{2} \mathbf{E}_{[i_{t-1}]} [\|x^* - x^{t-1}\|_2^2]. \end{aligned}$$

Dividing both sides of the above relation by  $\gamma_{t-1}\eta_{t-1}/2$ , we obtain

$$\begin{aligned} & \left( 1 + \frac{1-\gamma_{t-1}}{\gamma_{t-1}} \right) \mathbf{E}_{[i_t]} [\|x^* - x^t\|_2^2] + \frac{\tau_{t-1}}{2\eta_{t-1}\gamma_{t-1}} \mathbf{E}_{[i_t]} [\|y^* - y^t\|_2^2] \\ & \leq \|x^* - x^1\|_2^2 + \frac{\tau_{t-1}}{\eta_{t-1}\gamma_{t-1}} \|y^* - y^1\|_2^2 + \frac{1-\gamma_{t-1}}{\gamma_{t-1}} \mathbf{E}_{[i_{t-1}]} [\|x^* - x^{t-1}\|_2^2], \end{aligned} \quad (4.81)$$

which implies that

$$\left( 1 + \frac{1-\gamma_{t-1}}{\gamma_{t-1}} \right) \mathbf{E}_{[i_t]} [\|x^* - x^t\|_2^2] \leq \|x^* - x^1\|_2^2 + \frac{\tau_{t-1}}{\eta_{t-1}\gamma_{t-1}} \|y^* - y^1\|_2^2 + \frac{1-\gamma_{t-1}}{\gamma_{t-1}} \mathbf{E}_{[i_{t-1}]} [\|x^* - x^{t-1}\|_2^2].$$

For simplicity, let us denote

$$a = \frac{1-\gamma_{i-1}}{\gamma_{i-1}}, i = 1, \dots, t.$$

Then, using the previous inequality, and the definition of  $D$ , we have for any  $t \leq N$ ,

$$\begin{aligned} \mathbf{E}_{[i_t]} [\|x^* - x^t\|_2^2] & \leq \frac{D^2}{1+a} + \frac{a}{1+a} \mathbf{E}_{[i_{t-1}]} [\|x^* - x^{t-1}\|_2^2] \\ & \leq \frac{D^2}{1+a} + \frac{a}{1+a} \left( \frac{D^2}{1+a} + \frac{a}{1+a} \mathbf{E}_{[i_{t-2}]} [\|x^* - x^{t-2}\|_2^2] \right) \\ & = \frac{D^2}{1+a} \left( 1 + \frac{a}{1+a} \right) + \frac{a^2}{(1+a)^2} \mathbf{E}_{[i_{t-2}]} [\|x^* - x^{t-2}\|_2^2] \\ & \leq \dots \\ & \leq \frac{D^2}{1+a} \left( 1 + \frac{a}{1+a} + \dots + \frac{a^{t-2}}{(1+a)^{t-2}} \right) + \frac{a^{t-1}}{(1+a)^{t-1}} \|x^* - x^1\|_2^2 \\ & = \frac{D^2}{1+a} \frac{1 - \left(\frac{a}{1+a}\right)^{t-1}}{1 - \frac{a}{1+a}} + \frac{a^{t-1}}{(1+a)^{t-1}} \|x^* - x^1\|_2^2 \leq 2D^2. \end{aligned}$$

Plugging the above inequality to (4.81), for any  $t \leq N-1$ , we obtain

$$\frac{\tau_{t-1}}{2\eta_{t-1}\gamma_{t-1}} \mathbf{E}_{[i_t]} [\|y^* - y^t\|_2^2] \leq D^2 + \frac{1-\gamma_{t-1}}{\gamma_{t-1}} 2D^2,$$

which implies (4.75). ■

The following result provide an important bound on  $\mathbf{E}_{[i_N]} [Q_0(\hat{z}^N, z)]$  for the unbounded saddle point problems.

**Lemma 10** *Let  $z^t = (x^t, y^t), t = 1, 2, \dots, N$  be generated by the Algorithm 1 with  $x^1 = x^0$  and  $y^1 = \operatorname{argmax}_{y \in Y} \langle Ax^1, y \rangle - J(y)$ . Also assume that  $q_t, \tau_t, \eta_t$  and  $\gamma_t$  are set to (3.26)-(3.31) and (4.73). Then there exists a vector  $v_N$  such that*

$$\mathbf{E}_{[i_N]} [Q_0(\hat{z}^N, z) + \langle v_N, \hat{z}^N - z \rangle] \leq \left( \sum_{t=1}^{N-1} \gamma_t \right)^{-1} \left[ \frac{\gamma_{N-1}\eta_{N-1}}{2} \mathbf{E}_{[i_N]} [\|\hat{x}^N - x^1\|_2^2] + \frac{\tau_{N-1}}{2} \mathbf{E}_{[i_N]} [\|\hat{y}^N - y^1\|_2^2] \right]. \quad (4.82)$$

*Proof.* First note that

$$\begin{aligned}
\|x - x^1\|_2^2 - \|x - x^N\|_2^2 &= 2 \langle x^N - x^1, x \rangle + \|x^1\|_2^2 - \|x^N\|_2^2 \\
&= 2 \langle x^N - x^1, x - \hat{x}^N \rangle + 2 \langle x^N - x^1, \hat{x}^N \rangle + \|x^1\|_2^2 - \|x^N\|_2^2 \\
&= 2 \langle x^N - x^1, x - \hat{x}^N \rangle + \|x^1 - \hat{x}^N\|_2^2 - \|x^N - \hat{x}^N\|_2^2.
\end{aligned}$$

Using this identity in (3.37) and the fact that  $x^1 = x^0$  and  $y^1 = \operatorname{argmax}_{y \in Y} \langle Ax^1, y \rangle - J(y)$ , we obtain

$$\begin{aligned}
&\sum_{t=1}^{N-1} \gamma_t Q_0(z^{t+1}, z) + \langle Ax^N - Ax^{N-1}, \hat{y}^N - y \rangle \\
&+ \frac{\gamma_{N-1}\eta_{N-1}}{2} \langle x^N - x^1, \hat{x}^N - x \rangle + \frac{\tau_{N-1}}{2} \langle y^N - y^1, \hat{y}^N - y \rangle \\
\leq &\frac{\gamma_{N-1}\eta_{N-1}}{2} \|\hat{x}^N - x^1\|_2^2 - \frac{\gamma_{N-1}\eta_{N-1}}{2} \|\hat{x}^N - x^N\|_2^2 - \frac{\gamma_{N-1}\eta_{N-1}}{2} \|x^N - x^{N-1}\|_2^2 + \frac{\tau_{N-1}}{2} \|\hat{y}^N - y^1\|_2^2 \\
&- \frac{\tau_{N-1}}{2} \|\hat{y}^N - y^N\|_2^2 - \langle Ax^N - Ax^{N-1}, y^N - \hat{y}^N \rangle + \sum_{t=1}^{N-1} (\Delta'_t - \mathbf{E}_{i_t}[\Delta'_t]).
\end{aligned} \tag{4.83}$$

Denoting

$$v_N = \left( \sum_{t=1}^{N-1} \gamma_t \right)^{-1} \left( \frac{\gamma_{N-1}\eta_{N-1}}{2} (x^N - x^1), (Ax^N - Ax^{N-1}) + \frac{\tau_{N-1}}{2} (y^N - y^1) \right), \tag{4.84}$$

and using the fact that  $Q_0(z^{t+1}, z)$  is linear, we conclude from (4.83) that

$$\begin{aligned}
&\left( \sum_{t=1}^{N-1} \gamma_t \right) [Q_0(\hat{z}^N, z) + \langle v_N, \hat{z}^N - z \rangle] \\
\leq &\frac{\gamma_{N-1}\eta_{N-1}}{2} \|\hat{x}^N - x^1\|_2^2 - \frac{\gamma_{N-1}\eta_{N-1}}{2} \|\hat{x}^N - x^N\|_2^2 - \frac{\gamma_{N-1}\eta_{N-1}}{2} \|x^N - x^{N-1}\|_2^2 + \frac{\tau_{N-1}}{2} \|\hat{y}^N - y^1\|_2^2 \\
&- \frac{\tau_{N-1}}{2} \|\hat{y}^N - y^N\|_2^2 - \langle Ax^N - Ax^{N-1}, y^N - \hat{y}^N \rangle + \sum_{t=1}^{N-1} (\Delta'_t - \mathbf{E}_{i_t}[\Delta'_t]),
\end{aligned}$$

which together with the facts that

$$\begin{aligned}
-\langle Ax^N - Ax^{N-1}, y^N - \hat{y}^N \rangle &\leq \|A\|_2 \|x^N - x^{N-1}\|_2 \|y^N - \hat{y}^N\|_2 \\
&\leq \frac{\|A\|_2^2}{2\tau_{N-1}} \|x^N - x^{N-1}\|_2^2 + \frac{\tau_{N-1}}{2} \|y^N - \hat{y}^N\|_2^2,
\end{aligned}$$

and  $\gamma_{N-1}\eta_{N-1}\tau_{N-1} \geq \|A\|_2^2$ , then imply that

$$Q_0(\hat{z}^N, z) + \langle v_N, \hat{z}^N - z \rangle \leq \left( \sum_{t=1}^{N-1} \gamma_t \right)^{-1} \left[ \frac{\gamma_{N-1}\eta_{N-1}}{2} \|\hat{x}^N - x^1\|_2^2 + \frac{\tau_{N-1}}{2} \|\hat{y}^N - y^1\|_2^2 + \sum_{t=1}^{N-1} (\Delta'_t - \mathbf{E}_{i_t}[\Delta'_t]) \right].$$

The result now immediately follows by taking expectation w.r.t  $[i_N]$  on the both sides of the above inequality and noting that  $\mathbf{E}_{i_t}[\Delta'_t - \mathbf{E}_{i_t}[\Delta'_t]] = 0$ .  $\blacksquare$

The following theorem shows that the rate of convergence of the RPD algorithm for solving the unbounded saddle point problems.

**Theorem 11** *Let  $z^t = (x^t, y^t)$ ,  $t = 1, 2, \dots, N$  be generated by Algorithm 1 with  $x^1 = x^0$  and  $y^1 = \operatorname{argmax}_{y \in Y} \langle Ax^1, y \rangle - J(y)$ . Also assume that  $q_t, \tau_t, \eta_t$  and  $\gamma_t$  are set to (3.26)-(3.31) and (4.73).*

a) *For any  $N \geq 1$ , there exists a vector  $v_N$  such that*

$$\mathbf{E}_{[i_N]} [Q_0(\hat{z}^N, z) + \langle v_N, \hat{z}^N - z \rangle] \leq \frac{[3\gamma_{N-1}\eta_{N-1} + 2(2 - \gamma_{N-1})\eta_{N-1}]D^2}{\sum_{t=1}^{N-1} \gamma_t}, \tag{4.85}$$

$$\mathbf{E}_{[i_N]} [\|v_N\|_2] \leq \frac{KD}{\sum_{t=1}^{N-1} \gamma_t}, \tag{4.86}$$

where

$$K = 2\gamma_{N-1}\eta_{N-1} + 2\sqrt{\gamma_{N-1}\tau_{N-1}\eta_{N-1}} + \|A\|_2(1 + \sqrt{2}).$$

b) For any  $N \geq 1$ , there exists a vector  $\sigma(y)$  such that (4.86) holds,  $\mathbf{E}[\sigma(y)] = 0$  for any  $y \in Y$ , and

$$\mathbf{E}_{[i_N]} [\tilde{g}_{\sigma(y)}(\hat{z}^N, z, v_N)] \leq \frac{[3\gamma_{N-1}\eta_{N-1}+2(2-\gamma_{N-1})\eta_{N-1}]D^2}{\sum_{t=1}^{N-1} \gamma_t}. \quad (4.87)$$

*Proof.* We first show the part a). It follows from (4.74), (4.75), (4.76), and (4.77) that

$$\begin{aligned} & \mathbf{E}_{[i_N]} [\|(\gamma_{N-1}\eta_{N-1}(x^N - x^1), Ax^N - Ax^{N-1} + \tau_{N-1}(y^N - y^1))\|_2] \\ \leq & \mathbf{E}_{[i_N]} [\|Ax^N - Ax^{N-1}\|_2 + \gamma_{N-1}\eta_{N-1}\|x^N - x^1\|_2 + \tau_{N-1}\|y^N - y^1\|_2] \\ \leq & \mathbf{E}_{[i_N]} [\|A\|_2\|x^N - x^{N-1}\|_2 + \gamma_{N-1}\eta_{N-1}\|x^N - x^1\|_2 + \tau_{N-1}\|y^N - y^1\|_2] \\ \leq & \mathbf{E}_{[i_N]} [\gamma_{N-1}\eta_{N-1}(\|x^N - x^*\|_2 + \|x^1 - x^*\|_2) + \tau_{N-1}(\|y^N - y^*\|_2 + \|y^1 - y^*\|_2)] \\ & + \mathbf{E}_{[i_N]} [\|A\|_2(\|x^N - x^*\|_2 + \|x^{N-1} - x^*\|_2)] \\ \leq & 2\gamma_{N-1}\eta_{N-1}D + 2\sqrt{\gamma_{N-1}\tau_{N-1}\eta_{N-1}}D + \|A\|_2(D + \sqrt{2}D) = KD. \end{aligned}$$

The above inequality and the definition of  $v_N$  imply (4.86). On the other hand, using (4.74), (4.75), (4.76), and (4.77), we have

$$\begin{aligned} & \frac{\gamma_{N-1}\eta_{N-1}}{2}\mathbf{E}_{[i_N]} [\|\hat{x}^N - x^1\|_2^2] + \frac{\tau_{N-1}}{2}\mathbf{E}_{[i_N]} [\|\hat{y}^N - y^1\|_2^2] \\ \leq & \gamma_{N-1}\eta_{N-1}\mathbf{E}_{[i_N]} [\|\hat{x}^N - x^*\|_2^2 + \|x^* - x^1\|_2^2] + \tau_{N-1}\mathbf{E}_{[i_N]} [\|\hat{y}^N - y^*\|_2^2 + \|y^* - y^1\|_2^2] \\ = & \gamma_{N-1}\eta_{N-1}D^2 + \mathbf{E}_{[i_N]} [\gamma_{N-1}\eta_{N-1}\|\hat{x}^N - x^*\|_2^2 + \tau_{N-1}\|\hat{y}^N - y^*\|_2^2] \\ \leq & \gamma_{N-1}\eta_{N-1}D^2 + \frac{1}{\sum_{t=1}^{N-1} \gamma_t} \mathbf{E}_{[i_N]} \left[ \sum_{t=1}^{N-1} \gamma_t (\gamma_{t-1}\eta_{t-1}\|x^t - x^*\|_2^2 + \tau_{t-1}\|y^t - y^*\|_2^2) \right] \\ = & \gamma_{N-1}\eta_{N-1}D^2 + \frac{1}{\sum_{t=1}^{N-1} \gamma_t} \left[ \sum_{t=1}^{N-1} \gamma_t (\gamma_{t-1}\eta_{t-1}\mathbf{E}_{[i_t]} [\|x^t - x^*\|_2^2] + \tau_{t-1}\mathbf{E}_{[i_t]} [\|y^t - y^*\|_2^2]) \right] \\ \leq & \gamma_{N-1}\eta_{N-1}D^2 + \frac{1}{\sum_{t=1}^{N-1} \gamma_t} \left[ \sum_{t=1}^{N-1} \gamma_t \left( 2\gamma_{t-1}\eta_{t-1}D^2 + \tau_{t-1} \frac{2(2-\gamma_{t-1})\eta_{t-1}}{\tau_{t-1}} D^2 \right) \right] \\ = & (3\gamma_{N-1}\eta_{N-1} + 2(2-\gamma_{N-1})\eta_{N-1}) D^2. \end{aligned}$$

Using the above inequality in (4.82), we obtain (4.85). The proof of part b) is similar to that of Theorem 3.b) and hence the details are skipped.  $\blacksquare$

Below we specify a parameter setting that satisfies the assumptions in (3.26)-(3.31) and (4.73) and leads to an optimal rate of convergence for the RPD algorithm in terms of its dependence on  $N$ .

**Corollary 12** Let  $z^t = (x^t, y^t)$ ,  $t = 1, 2, \dots, N$  be generated by Algorithm 1 with  $x^1 = x^0$  and  $y^1 = \operatorname{argmax}_{y \in Y} \langle Ax^1, y \rangle - J(y)$ . Also assume that  $\gamma_t, q_t, \tau_t$  and  $\eta_t$  are set to

$$q_t = p, \quad \forall t \leq N, \quad (4.88)$$

$$\gamma_t = \frac{1}{p}, \quad \forall t = 1, 2, \dots, N-2, \quad \text{and} \quad \gamma_{N-1} = 1, \quad (4.89)$$

$$\tau_t = \|A\|_2 p^{3/2}, \quad i = 1, \dots, N-1, \quad (4.90)$$

$$\eta_t = \|A\|_2 p^{3/2}, \quad i = 1, \dots, N-2, \quad \text{and} \quad \eta_{N-1} = \|A\|_2 p^{1/2}. \quad (4.91)$$

Then for any  $N \geq 1$ , there exists a vector  $v_N$  such that

$$\mathbf{E}_{[i_N]} [Q_0(\hat{z}^N, z) + \langle v_N, \hat{z}^N - z \rangle] \leq \frac{5p^{3/2}\|A\|_2 D^2}{N+p-2}, \quad (4.92)$$

$$\mathbf{E}_{[i_N]} [\|v_N\|_2] \leq \frac{p}{N+p-2} \left( 4p^{1/2} + (1 + \sqrt{2}) \right) \|A\|_2 D. \quad (4.93)$$

Moreover, for any  $N \geq 1$ , there exists a vector  $\sigma(y)$  such that (4.93) holds,  $\mathbf{E}[\sigma(y)] = 0$  for any  $y \in Y$ , and

$$\mathbf{E}_{[i_N]} [\tilde{g}_{\sigma(y)}(\hat{z}^N, z, v_N)] \leq \frac{5p^{3/2}\|A\|_2 D^2}{N+p-2}. \quad (4.94)$$

*Proof.* It is easy to verify that  $\gamma_t, q_t, \tau_t$  and  $\eta_t$  defined in (4.88)-(4.91) satisfy (3.26)-(3.31) and (4.73). We also have

$$\sum_{t=1}^{N-1} \gamma_t = \frac{(N+p-2)}{p}.$$

Plugging this identity into (4.85)-(4.87), we obtain (4.92)-(4.94) respectively.  $\blacksquare$

A few remarks about the results obtained in Theorem 11 and Corollary 12 are in place. First, in the view of (4.92), the total number of iterations required by the RPD algorithm to find an  $\epsilon$ -solution of problem (1.1), i.e., a point  $\hat{z} \in Z$  such that  $\mathbf{E}[Q_0(\hat{z}, z) + \langle v_N, \hat{z}^N - z \rangle] \leq \epsilon$  for any  $z \in Z$ , can be bounded by  $\mathcal{O}(p^{3/2}\|A\|_2\Omega_X\Omega_Y/\epsilon)$ . Second, similar to the bounded problems, these results are new and optimal in terms of its dependence on  $\epsilon$  for a given  $p$  (see discussions in [9]). To the best of our knowledge, this is the first time such an optimal rate of convergence is obtained in the literature for a randomized algorithm for solving the saddle point problem (1.1)-(1.2) with unbounded domains.

## 4.2 Non-Euclidean randomized primal-dual methods

In this subsection, we show that by replacing the usual Euclidean distance by generalized non-Euclidean prox-functions, Algorithm 1 can be adaptive to different geometry of the feasible sets

Recall that a function  $\omega_i : Y_i \rightarrow R$  is a distance generating function [37] with modulus  $\alpha_i$  with respect to  $\|\cdot\|_i$ , if  $\omega_i$  is continuously differentiable and strongly convex with parameter  $\alpha_i$  with respect to  $\|\cdot\|_i$ . Without loss of generality, we assume that  $\alpha_i = 1$  for any  $i = 1, \dots, b$ , because we can always rescale  $\omega_i(y)$  to  $\bar{\omega}_i(y) = \omega_i(y)/\alpha_i$  in case  $\alpha_i \neq 1$ . Therefore, we have

$$\langle y - z, \nabla\omega_i(y) - \nabla\omega_i(z) \rangle \geq \|y - z\|_i^2 \quad \forall y, z \in Y_i.$$

The prox-function associated with  $\omega_i$  is given by

$$V_i(z, y) = \omega_i(y) - [\omega_i(z) + \langle \nabla\omega_i(z), y - z \rangle] \quad \forall y, z \in Y_i. \quad (4.95)$$

The prox-function  $V_i(\cdot, \cdot)$  is also called the Bregman's distance, which was initially studied by Bregman [4]. Suppose that the set  $Y_i$  is bounded, the distance generating function  $\omega_i$  also gives rise to the diameter of  $Y_i$ , which will be used frequently in our convergence analysis:

$$\mathcal{D}_{\omega_i, Y_i} := \max_{y \in Y_i} \omega_i(y) - \min_{y \in Y_i} \omega_i(y). \quad (4.96)$$

For the sake of notational convenience, sometimes we simply denote  $\mathcal{D}_{\omega_i, Y_i}$  by  $\mathcal{D}_i$ ,  $V(y, z) = \sum_{i=1}^p V_i(y^{(i)}, z^{(i)})$ ,  $\forall y, z \in Y$ , and  $D_Y = \sum_{i=1}^p \mathcal{D}_i$ . Let  $y_1^{(i)} = \operatorname{argmin}_{y \in Y_i} \omega_i(y)$ ,  $i = 1, \dots, b$ . We can easily see that for any  $y \in Y$ ,

$$\begin{aligned} \|y_1^{(i)} - y^{(i)}\|_i^2/2 &\leq V_i(y_1^{(i)}, y^{(i)}) = \omega_i(y^{(i)}) - \omega_i(y_1^{(i)}) - \langle \nabla\omega_i(y_1^{(i)}), y^{(i)} - y_1^{(i)} \rangle \\ &\leq \omega_i(y^{(i)}) - \omega_i(y_1^{(i)}) \leq \mathcal{D}_i. \end{aligned} \quad (4.97)$$

Moreover, we define  $\|y\|^2 = \|y^{(1)}\|_1^2 + \dots + \|y^{(p)}\|_p^2$  and denote its conjugate by  $\|y\|_*^2 = \|y^{(1)}\|_{1,*}^2 + \dots + \|y^{(p)}\|_{p,*}^2$ . Similarly, letting  $\omega : X \rightarrow \mathbb{R}$  be continuously differentiable and strongly convex w.r.t  $\|\cdot\|$  with modulus 1, we define the prox-function  $V(\cdot, \cdot)$  associated with  $\omega$  and use  $D_X$  to denote the diameter of  $X$ .

We are now ready to describe a non-Euclidean variant of Algorithm 1, which is obtained by replacing the Euclidean distances used in the two subproblems (3.10) and (3.11) in Step 2 of Algorithm 1 with the Bregman's distances in (4.98) and (4.99).

---

**Algorithm 2** The non-Euclidean RPD Method

---

Let  $z^1 = (x^1, y^1) \in X \times Y$ , and nonnegative stepsizes  $\{\tau_t\}$ ,  $\{\eta_t\}$ , parameters  $\{q_t\}$ , and weights  $\{\gamma_t\}$  be given. Set  $\bar{x}^1 = x^1$ .

**for**  $t = 1, \dots, N$  **do**

1. Generate a random variable  $i_t$  uniformly from  $\{1, 2, \dots, p\}$ .
2. Update  $y^{t+1}$  and  $x^{t+1}$  by

$$y_i^{t+1} = \begin{cases} \operatorname{argmin}_{y_{i_t} \in Y_{i_t}} \langle -U_{i_t} A \bar{x}^t, y \rangle + J_{i_t}(y_{i_t}) + \tau_t V_{i_t}(y_{i_t}, y_{i_t}^t), & i = i_t, \\ y_i^t, & i \neq i_t. \end{cases} \quad (4.98)$$

$$x^{t+1} = \operatorname{argmin}_{x \in X} h(x) + \langle x, A^T y^{t+1} \rangle + \eta_t V(x, x^t). \quad (4.99)$$

$$\bar{x}^{t+1} = q_t(x^{t+1} - x^t) + x^{t+1}. \quad (4.100)$$

**end for**

**Output:** Set

$$\hat{z}^N = \left( \sum_{t=1}^{N-1} \gamma_t \right)^{-1} \sum_{t=1}^{N-1} \gamma_t z^{t+1}. \quad (4.101)$$


---

We will show that the non-Euclidean RPD algorithms exhibit similar convergence properties to the Euclidean RPD algorithm for solving general bilinear saddle point problems with bounded feasible sets, but they can be more flexible on the selection of the norms and distance generating functions.

First, the following result generalizes Proposition 1.

**Proposition 13** *Let  $\{y^t\}_{t \geq 1}$  and  $\{x^t\}_{t \geq 1}$  be generated by Algorithm 2. Then for any  $z \in Z$ , we have*

$$\begin{aligned} & \gamma_t Q_0(z^{t+1}, z) + \langle \gamma_t A x^{t+1} - A x^t, y^{t+1} - y \rangle + (\gamma_t - 1) [J(y) - J(y^{t+1})] - \Delta_t \\ & \leq \gamma_t \eta_t [V(x, x^t) - V(x^t, x^{t+1}) - V(x, x^{t+1})] + \tau_t [V(y, y^t) - V(y, y^{t+1}) - V(y^t, y^{t+1})], \end{aligned} \quad (4.102)$$

where  $\Delta_t$  is defined in (3.16).

*Proof.* By the optimality condition of problem (4.99), for all  $x \in X$ , we have

$$h(x^{t+1}) - h(x) + \langle x^{t+1} - x, A^T y^{t+1} \rangle + \eta_t V(x^t, x^{t+1}) + \eta_t V(x, x^{t+1}) \leq \eta_t V(x, x^t). \quad (4.103)$$

Similarly, by the optimality condition of problem (4.98), for all  $y \in Y$ , we have

$$\langle -U_{i_t} A \bar{x}^t, y^{t+1} - y \rangle + J_{i_t}(y_{i_t}^{t+1}) - J_{i_t}(y_{i_t}) + \tau_t V_{i_t}(y_{i_t}^t, y_{i_t}^{t+1}) + \tau_t V_{i_t}(y_{i_t}, y_{i_t}^{t+1}) \leq \tau_t V_{i_t}(y_{i_t}, y_{i_t}^t). \quad (4.104)$$

The result follows by using an argument similar to the one used in the proof of Proposition 1 by replacing the Euclidean distances with Bregman's distances and noting that

$$\begin{aligned} V_{i_t}(y_{i_t}^t, y_{i_t}^{t+1}) &= V(y^t, y^{t+1}), \\ V_{i_t}y_{i_t}, y_{i_t}^t - V_{i_t}(y_{i_t}, y_{i_t}^{t+1}) &= V(y, y^t) - V(y, y^{t+1}). \end{aligned}$$

■

The following lemma provides an upper bound of  $\mathbf{E}_{i_t}[\Delta_t]$ .

**Lemma 14** *If  $i_t$  is uniformly distributed on  $\{1, 2, \dots, p\}$ , then*

$$\begin{aligned} \mathbf{E}_{i_t}[\Delta_t] \leq & \left\langle \left( \frac{1}{p}q_{t-1} - \frac{p-1}{p} \right) Ax^t - \frac{1}{p}q_{t-1}Ax^{t-1}, y^t - y \right\rangle + \frac{p-1}{p} [J(y^t) - J(y)] \\ & + \frac{q_{t-1}^2 \|A\|^2}{2p\tau_t} \|x^t - x^{t-1}\|^2 + \frac{\tau_t}{2} \mathbf{E}_{i_t} [\|y^{t+1} - y^t\|^2]. \end{aligned}$$

*Proof.* The proof is similar to the one of Lemma 2 except that we now replace the Euclidean distances by Bregman's distances, and that in (3.25), we use the fact that  $\|x - z\|^2/2 \leq V(x, z)$ , i.e.,

$$\begin{aligned} \mathbf{E}_{i_t} [\langle q_{t-1}U_{i_t}A(x^t - x^{t-1}), y^{t+1} - y^t \rangle] &\leq \mathbf{E}_{i_t} [q_{t-1}\|U_{i_t}A(x^t - x^{t-1})\| \|y^{t+1} - y^t\|] \\ &\leq \mathbf{E}_{i_t} \left[ \frac{q_{t-1}^2}{2\tau_t} \|U_{i_t}A(x^t - x^{t-1})\|^2 + \frac{\tau_t}{2} \|y^{t+1} - y^t\|^2 \right] \\ &\leq \frac{q_{t-1}^2}{2\tau_t} \mathbf{E}_{i_t} [\|U_{i_t}A(x^t - x^{t-1})\|^2] + \frac{\tau_t}{2} \mathbf{E}_{i_t} [\|y^{t+1} - y^t\|^2] \\ &= \frac{q_{t-1}^2}{2p\tau_t} \|A(x^t - x^{t-1})\|^2 + \frac{\tau_t}{2} \mathbf{E}_{i_t} [\|y^{t+1} - y^t\|^2] \\ &\leq \frac{q_{t-1}^2 \|A\|^2}{2p\tau_t} \|x^t - x^{t-1}\|^2 + \frac{\tau_t}{2} \mathbf{E}_{i_t} [\|y^{t+1} - y^t\|^2] \\ &\leq \frac{q_{t-1}^2 \|A\|^2}{p\tau_t} V(x^t, x^{t-1}) + \tau_t \mathbf{E}_{i_t} [V(y^{t+1}, y^t)]. \end{aligned}$$

■

Theorem 15 below describes some convergence properties of the non-Euclidean RPD methods.

**Theorem 15** *Suppose that the starting point  $z^1$  is chosen such that  $x^1 = x^0$  and  $y^1 = \operatorname{argmax}_{y \in Y} \langle Ax^1, y \rangle - J(y)$ . Also assume that the parameters  $q_t, \gamma_t, \tau_t$  and  $\eta_t$  are set to (3.26)-(3.31). Then, for any  $N \geq 1$ , we have*

$$\mathbf{E}_{[i_N]}[Q_0(\hat{z}^N, z)] \leq \left( \sum_{t=1}^{N-1} \gamma_t \right)^{-1} \left[ \frac{\gamma_1 \eta_1}{2} D_X + \frac{\tau_1}{2} D_Y \right], \quad \forall z \in Z. \quad (4.105)$$

where  $\hat{z}^N$  is defined in (3.13) and the expectation is taken w.r.t. to  $i_{[N]} = (i_1, \dots, i_{N-1})$ .

*Proof.* The proof is almost identical to that of Theorem 3 except that we replace the Euclidean distances with Bregman's distances and that in (3.38) we use the fact  $V(x, x_1) \leq D_X$  and  $V(y, y_1) \leq D_Y$ , i.e.,

$$\begin{aligned} -\langle Ax^N - Ax^{N-1}, y^N - y \rangle &\leq \frac{\|A\|^2}{2\tau_{N-1}} \|x^N - x^{N-1}\|^2 + \frac{\tau_{N-1}}{2} \|y^N - y\|^2 \\ &\leq \frac{\|A\|^2}{\tau_{N-1}} V(x^N, x^{N-1}) + \tau_{N-1} V(y^N, y). \end{aligned}$$

■

It should be noted that we can also establish the convergence of the generalized algorithm for smooth bilinear saddle point problems. However, it is still not clear to us whether Algorithm 2 can be generalized to the case when neither  $X$  nor  $Y$  are bounded.

## 5 RPD for linearly constrained problems and its relation to ADMM

Our goal of this section is to show that Algorithm 1 applied to the linearly constrained optimization problems can be viewed exactly as a randomized proximal alternating direction of multiplier method (ADMM).

More specifically, consider the following optimization problem

$$\begin{aligned} \min \quad & f_1(x_1) + f_2(x_2) + \dots + f_p(x_p) \\ \text{s.t.} \quad & A_1x_1 + A_2x_2 + \dots + A_px_p = b \end{aligned} \tag{5.106}$$

Chen et. al. show in [7] that a direct extension of ADMM does not necessarily converge for solving this type of problem whenever  $p \geq 3$ . More precisely, for the case  $p \geq 3$ , it is required that the given coefficient matrices  $A_i$  satisfy some orthogonality assumptions. In [23], Luo and Hong proposed a variant of ADMM, namely the proximal ADMM, and proved its asymptotical convergence under the strong convexity assumptions about the objective function. However, to the best of our knowledge, there does not exist a proof for the convergence for the proximal ADMM method when the strong convexity assumption is removed. On the other hand, the RPD method, which will be shown to be equivalent to a randomized version of the proximal ADMM method, exhibits an  $\mathcal{O}(1/N)$  rate of convergence for solving problem (5.106) without requiring any assumptions about the matrices  $A_1, A_2, \dots, A_p$ , as well as the strong convexity assumptions about  $f_i, i = 1, \dots, p$ .

Let us first formally state these two algorithms. Observing that problem (5.106) is equivalent to

$$\min_{y \in Y} \max_{x \in X} \{ \langle y, b \rangle - \langle y, \sum_{i=1}^p A_i x_i \rangle - \sum_{i=1}^p f_i(x_i) \}, \tag{5.107}$$

where  $Y = \mathbb{R}^m$ , we can specialize Algorithm 1 applied to problem (5.107) as shown Algorithm 3. On the other hand, noting that the augmented Lagrangian function of (5.106) is given by

$$L(x, y) = \left\{ \sum_{i=1}^p f_i(x_i) + \langle y, \sum_{i=1}^p A_i x_i - b \rangle + \frac{\rho}{2} \left\| \sum_{i=1}^p A_i x_i - b \right\|^2 \right\}, \tag{5.108}$$

we can state the proximal ADMM method for solving problem (5.106) as shown in Algorithm 4. It is easy to see that (5.111) can be rewritten as

$$y^{t+1} = y^t + \frac{1}{\tau_t} \left( \sum_{i=1}^p A_i x_i^{t+1} - b \right),$$

which implies that

$$\bar{y}^{t+1} = y^{t+1} + \frac{q_t}{\tau_t} \left( \sum_{i=1}^p A_i x_i - b \right). \tag{5.109}$$

In view of (5.109), if only a randomly selected block  $x_{i_t}^{t+1}$  is updated in the Step 2 of the proximal ADMM method instead of all blocks of  $x^{t+1}$ , then the randomized version of (5.113) and (5.110) are equivalent in case  $\rho = \frac{q_t - 1}{\tau_t - 1}$ . Therefore, we conclude that Algorithm 1 applied to problem (5.107) is equivalent to a randomized version of the proximal ADMM method for solving linearly constrained problems (5.106).



---

**Algorithm 4** Proximal alternating direction of multiplier methods

---

Let  $z = (x^1, y^1) \in X \times Y$  and stepsizes  $\{\eta_t\}_{t \geq 1}$ .

**for**  $t = 1, \dots, N$  **do**

Update  $y^{t+1}$  and  $x^{t+1}$  by

$$x_i^{t+1} = \arg \min_{x_i \in X_i} f_i(x_i) + \langle y_t, A_i x_i \rangle + \rho \langle \sum_{j < i} A_j x_j^{t+1} + \sum_{j \geq i} A_j x_j^t - b, A_i x_i \rangle + \frac{\eta_t}{2} \|x_i - x_i^t\|_2^2, i = 1, \dots, p. \quad (5.113)$$

$$y^{t+1} = y^t + \rho \left( \sum_{i=1}^p A_i x_i^{t+1} - b \right). \quad (5.114)$$

**end for**

---

---

**Algorithm 3** Randomized primal-dual Methods for problem (5.107)

---

Let  $z^1 = (x^1, y^1) \in X \times Y$  and stepsizes  $\{\gamma_t\}_{t \geq 1}$ ,  $\{q_t\}_{t \geq 1}$ ,  $\{\tau_t\}_{t \geq 1}$ ,  $\{\eta_t\}_{t \geq 1}$ .

**for**  $t = 1, \dots, N$  **do**

1. Generate a random variable  $i_t$  from  $\{1, \dots, p\}$ .
2. Update  $y^{t+1}$  and  $x^{t+1}$  by

$$x_i^{t+1} = \begin{cases} \operatorname{argmin}_{x_{i_t} \in X_{i_t}} f_{i_t}(x_{i_t}) + \langle \bar{y}_t, A_{i_t} x_{i_t} \rangle + \frac{\eta_t}{2} \|x_{i_t} - x_{i_t}^t\|_2^2, & i = i_t, \\ x_i^t, & i \neq i_t. \end{cases} \quad (5.110)$$

$$y^{t+1} = \operatorname{argmin}_{y \in Y} \langle y, b \rangle - \langle y, \sum_{i=1}^p A_i x_i^{t+1} \rangle + \frac{\tau_t}{2} \|y - y^t\|_2^2. \quad (5.111)$$

$$\bar{y}^{t+1} = q_t(y^{t+1} - y^t) + y^{t+1}. \quad (5.112)$$

**end for**

---

In order to understand its practical performance for solving the worst-case instances in [7], we implement Algorithm 1 for solving the linearly constrained problem (5.106) with  $b = 0$  and  $f_i(x_i) = 0$ ,  $i = 1, 2, \dots, p$ . Moreover, we assume that  $A_i, i = 1, \dots, p$  are set to  $A_1 = (1; 1; \dots; 1), A_2 = (1; \dots; 1; 2), \dots, A_p = (1; 2; 2; \dots; 2)$ . Under the above settings, problem (5.106) is equivalent to a homogenous linear system with  $p$  variables

$$\sum_{i=1}^p A_i x_i = 0, \quad (5.115)$$

where  $A_i, i = 1, 2, \dots, p$  are nonsingular. Problem (5.115) has a unique solution  $x^* = (0; 0; \dots; 0) \in \mathbb{R}^n$ . The problem constructed above slightly generalizes the counter example in [7]. As shown in Table 1, while the original ADMM does not necessarily converge in solving the above problem even with  $p = 3$ , Algorithm 1 converges to the optimal solution  $x^*$  for all different values of  $p$  that we have tested.

## 6 Concluding remarks

In this paper, we present a new randomized algorithm, namely the randomized primal-dual method, for solving a class of bilinear saddle point problems. Each iteration of the RPD method requires

Table 1: Results of Algorithm 1 for solving problem (5.115)

p	$\ x^{100} - x^*\ $	$\ x^{1,000} - x^*\ $	$\ x^{10,000} - x^*\ $	$\ x^{100,000} - x^*\ $
10	2.0608	1.1416	0.2674	0.0396
20	4.2308	1.1438	1.6588	0.4711
50	7.0277	6.6469	2.2886	2.1143

to solve only one subproblem rather than all subproblems as in the original primal-dual algorithms. The RPD method does not require strong convexity assumptions about the objective function and/or boundedness assumptions about the feasible sets. Moreover, based on a new primal-dual termination criterion, we show that this algorithm exhibits an  $\mathcal{O}(1/N)$  rate of convergence for both bounded and unbounded saddle point problems and an  $\mathcal{O}(1/N^2)$  rate of convergence for smooth saddle point problems. Extension for the non-Euclidean setting and the relation to the ADMM method have also been discussed in this paper.

It is worth noting that there exist a few possible extensions of this work. Firstly, for the case when  $h(x)$  is not necessarily simple, but a general smooth convex, one can modify (3.11) in Step 2 of Algorithm 1 by replacing  $h(x)$  with its linear approximation as suggested in [9]. Secondly, this paper focuses on the case when the dual space has multiple blocks. However, it is possible to apply block decomposition for both the primal and dual spaces whenever the feasible sets  $X$  and  $Y$  are decomposable. Finally, it will be interesting to see if the rate of convergence for the RPD methods can be further improved by using non-uniform distribution for the random variables  $i_t$ .

## References

- [1] A. Beck and L. Tetrushvili. On the convergence of block coordinate descent type methods. Technical report. submitted to *SIAM Journal on Optimization*.
- [2] D. Blatt, A. Hero, and H. Gauchman. A convergent incremental gradient method with a constant step size. *SIAM Journal on Optimization*, 18(1):29–51, 2007.
- [3] Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, and Jonathan Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine Learning*, 3(1):1–122, 2011.
- [4] L.M. Bregman. The relaxation method of finding the common point convex sets and its application to the solution of problems in convex programming. *USSR Comput. Math. Phys.*, 7:200–217, 1967.
- [5] Regina S Burachik, Alfredo N Iusem, and Benar Fux Svaiter. Enlargement of monotone operators with applications to variational inequalities. *Set-Valued Analysis*, 5(2):159–180, 1997.
- [6] A. Chambolle and T. Pock. A first-order primal-dual algorithm for convex problems with applications to imaging. *J. Math. Imaging Vision*, 40:120–145, 2011.
- [7] Caihua Chen, Bingsheng He, Yinyu Ye, and Xiaoming Yuan. The direct extension of admm for multi-block convex minimization problems is not necessarily convergent. *Optimization Online*, 2013.

- [8] Y. Chen, G. Lan, and Y. Ouyang. Accelerated schemes for a class of variational inequalities. *Mathematical Programm, Series B*, 2014. submitted.
- [9] 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.
- [10] C. D. Dang and G. Lan. Stochastic block mirror descent methods for nonsmooth and stochastic optimization. *SIAM Journal on Optimization*, 2015. to appear.
- [11] Jr. Douglas, Jim and Jr. Rachford, H. H. On the numerical solution of heat conduction problems in two and three space variables. *Transactions of the American Mathematical Society*, 82(2):pp. 421–439, 1956.
- [12] Jonathan Eckstein and Dimitri P. Bertsekas. On the douglas-rachford splitting method and the proximal point algorithm for maximal monotone operators. *Mathematical Programming*, 55(1-3):293–318, 1992.
- [13] E. Esser, X. Zhang, and T.F. Chan. A general framework for a class of first order primal-dual algorithms for convex optimization in imaging science. *SIAM Journal on Imaging Sciences*, 3(4):1015–1046, 2010.
- [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.
- [15] Daniel Gabay. Chapter ix applications of the method of multipliers to variational inequalities. *Studies in mathematics and its applications*, 15:299–331, 1983.
- [16] S. Ghadimi and G. Lan. Optimal stochastic approximation algorithms for strongly convex stochastic composite optimization, I: a generic algorithmic framework. *SIAM Journal on Optimization*, 22:1469–1492, 2012.
- [17] Roland Glowinski and A Marroco. Sur l’approximation, par éléments 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-Modélisation Mathématique et Analyse Numérique*, 9(R2):41–76, 1975.
- [18] D. Goldfarb and S. Ma. Fast multiple-splitting algorithms for convex optimization. *SIAM Journal on Optimization*, 22(2):533–556, 2012.
- [19] Donald Goldfarb, Shiqian Ma, and Katya Scheinberg. Fast alternating linearization methods for minimizing the sum of two convex functions. *Mathematical Programming*, 141(1-2):349–382, 2013.
- [20] B. He and X. Yuan. On the  $\mathcal{O}(1/n)$  convergence rate of the douglasrachford alternating direction method. *SIAM Journal on Numerical Analysis*, 50(2):700–709, 2012.
- [21] Bingsheng He and Xiaoming Yuan. On the  $\mathcal{o}(1/n)$  convergence rate of the douglas-rachford alternating direction method. *SIAM Journal on Numerical Analysis*, 50(2):700–709, 2012.

- [22] N. He, A. Juditsky, and A. Nemirovski. Mirror prox algorithm for multi-term composite minimization and semi-separable problems. *Computational Optimization and Applications*, 2014. submitted.
- [23] M. Hong and Z.-Q. Luo. On the Linear Convergence of the Alternating Direction Method of Multipliers. *ArXiv e-prints*, August 2012.
- [24] G. Lan. An optimal method for stochastic composite optimization. *Mathematical Programming*, 133(1):365–397, 2012.
- [25] G. Lan, Z. Lu, and R. D. C. Monteiro. Primal-dual first-order methods with  $\mathcal{O}(1/\epsilon)$  iteration-complexity for cone programming. *Mathematical Programming*, 126:1–29, 2011.
- [26] G. Lan and R. D. C. Monteiro. Iteration-complexity of first-order augmented lagrangian methods for convex programming. *Mathematical Programming*.
- [27] G. Lan and R. D. C. Monteiro. Iteration-complexity of first-order penalty methods for convex programming. *Mathematical Programming*, 138:115–139, 2013.
- [28] D. Leventhal and A. S. Lewis. Randomized methods for linear constraints: Convergence rates and conditioning. *Mathematics of Operations Research*, 35:641–654, 2010.
- [29] P. L. Lions and B. Mercier. Splitting algorithms for the sum of two nonlinear operators. *SIAM Journal on Numerical Analysis*, 16(6):pp. 964–979, 1979.
- [30] Z. Lu and L. Xiao. On the complexity analysis of randomized block-coordinate descent methods. Manuscript, 2013.
- [31] Mehrdad Mahdavi and Rong Jin. Mixedgrad: An  $O(1/T)$  convergence rate algorithm for stochastic smooth optimization. *CoRR*, abs/1307.7192, 2013.
- [32] R.D.C. Monteiro and B.F. Svaiter. On the complexity of the hybrid proximal extragradient method for the iterates and the ergodic mean. Manuscript, School of ISyE, Georgia Tech, Atlanta, GA, 30332, USA, March 2009.
- [33] R.D.C. Monteiro and B.F. Svaiter. Complexity of variants of tsengs modified f-b splitting and korpelevich’s methods for hemi-variational inequalities with applications to saddle-point and convex optimization problems. Manuscript, School of ISyE, Georgia Tech, Atlanta, GA, 30332, USA, June 2010.
- [34] R.D.C. Monteiro and B.F. Svaiter. On the complexity of the hybrid proximal projection method for the iterates and the ergodic mean. *SIAM Journal on Optimization*, 20:2755–2787, 2010.
- [35] Renato DC Monteiro and Benar F Svaiter. Iteration-complexity of block-decomposition algorithms and the alternating direction method of multipliers. *SIAM Journal on Optimization*, 23(1):475–507, 2013.
- [36] A. S. Nemirovski. Prox-method with rate of convergence  $o(1/t)$  for variational inequalities with lipschitz continuous monotone operators and smooth convex-concave saddle point problems. *SIAM Journal on Optimization*, 15:229–251, 2005.

- [37] A. S. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization*, 19:1574–1609, 2009.
- [38] Y. E. Nesterov. Smooth minimization of nonsmooth functions. *Mathematical Programming*, 103:127–152, 2005.
- [39] Y. E. Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. Technical report, Center for Operations Research and Econometrics (CORE), Catholic University of Louvain, February 2010.
- [40] Y. E. Nesterov. Subgradient methods for huge-scale optimization problems. Technical report, Center for Operations Research and Econometrics (CORE), Catholic University of Louvain, February 2012.
- [41] Y. Ouyang, Y. Chen, G. Lan, and E. Pasiliao. An accelerated linearized alternating direction method of multipliers. *SIAM Journal on Imaging Sciences*, 2014. to appear.
- [42] P. Richtárik and M. Takáč. Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function. *Mathematical Programming*, 2012. to appear.
- [43] R. Rockafellar. Monotone operators and the proximal point algorithm. *SIAM Journal on Control and Optimization*, 14(5):877–898, 1976.
- [44] M. Schmidt, N. L. Roux, and F. Bach. Minimizing finite sums with the stochastic average gradient. Technical report, September 2013.
- [45] S. Shalev-Shwartz and T. Zhang. Accelerated proximal stochastic dual coordinate ascent for regularized loss minimization. *Mathematical Programming*, 2015. to appear.
- [46] Shai Shalev-Shwartz and Tong Zhang. Stochastic dual coordinate ascent methods for regularized loss. *J. Mach. Learn. Res.*, 14(1):567–599, February 2013.
- [47] T. Suzuki. Stochastic Dual Coordinate Ascent with Alternating Direction Multiplier Method. *ArXiv e-prints*, November 2013.
- [48] Shuzhong Zhang Tianyi Lin, Shiqian Ma. On the Global Linear Convergence of the ADMM with Multi-Block Variables. *ArXiv e-prints*, 2014.
- [49] Yuchen Zhang and Lin Xiao. Stochastic primal-dual coordinate method for regularized empirical risk minimization. Manuscript, September 2014.