

Separable Approximations and Decomposition Methods for the Augmented Lagrangian

Rachael Tappenden Peter Richtárik Burak Büke *

Abstract

In this paper we study decomposition methods based on separable approximations for minimizing the augmented Lagrangian. In particular, we study and compare the Diagonal Quadratic Approximation Method (DQAM) of Mulvey and Ruszczyński [20] and the Parallel Coordinate Descent Method (PCDM) of Richtárik and Takáč [30]. We show that the two methods are equivalent for feasibility problems up to the selection of a step-size parameter. Further, we prove an improved complexity bound for PCDM under strong convexity, and show that this bound is at least $8(L'/\bar{L})(\omega - 1)^2$ times better than the best known bound for DQAM, where ω is the degree of partial separability and L' and \bar{L} are the maximum and average of the block Lipschitz constants of the gradient of the quadratic penalty appearing in the augmented Lagrangian.

1 Introduction

With the rise and ubiquity of digital and data technology, practitioners in nearly all industries need to solve optimization problems of increasingly larger sizes. As a consequence, new tools and methods are required to solve these big data problems, and to do so efficiently.

In this work, we are concerned with convex optimization problems with an objective function that is separable into blocks of variables and where these blocks are linked by a subset of constraints which nevertheless make the problem nonseparable. Nonseparability is a source of difficulty in solving these very large optimization problems. This structure is particularly relevant in stochastic optimization problems where each block relates to a certain scenario and involves only variables related to that particular scenario. The objective function expressed as an expectation is separable in these blocks and the linking constraints (called non-anticipativity constraints) encode the natural requirement that decisions be based only on information available at the time of decision making. Applications that can be modeled as large scale stochastic optimization problems include multicommodity network flow problems, financial planning problems and airline routing.

A classical approach to solving such problems is to use the augmented Lagrangian by relaxing the linking constraints. The augmented Lagrangian idea was first introduced independently by Hestenes [13] and Powell [26] and convergence of the associated augmented Lagrangian method was established later by Rockafellar [32, 33]. The method remains popular today with studies on, for example, convergence of the method [4] and applying the method to design optimization [11]. Advantages of the augmented Lagrangian method approach include the simplicity and stability of the multiplier iterations, the possibility of starting from an arbitrary multiplier, and the fact that

*All authors: James Clerk Maxwell Building, School of Mathematics, The University of Edinburgh, United Kingdom. The work of all three authors was supported by the EPSRC grant EP/I017127/1 (Mathematics for Vast Digital Resources). The work of PR and RT was also partially supported by the Centre for Numerical Algorithms and Intelligent Software (funded by EPSRC grant EP/G036136/1 and the Scottish Funding Council).

there is no primal master problem to solve. However, the augmented Lagrangian is nonseparable, so the problem is still difficult to solve.

The nonseparability of the augmented Lagrangian has motivated the development of decomposition techniques. In an early work, Stephanopoulos and Westerberg [40] suggest decomposing the augmented Lagrangian using linear approximations and Watanabe et al. [45] use a transformation method to deal with the nonseparable cross products. In a more recent line of work, Ruszczyński [35, 36] and Mulvey and Ruszczyński [20, 21] propose and analyze a diagonal quadratic approximation (DQA) to the augmented Lagrangian and an associated diagonal quadratic approximation method (DQAM). By approximating the original problem by one that is separable into blocks, these techniques make a significant difference in terms of solvability because the problem is broken into a number of problems of a more manageable size. Decomposition techniques have become even more attractive with the advances in parallel computing: since the decomposed subproblems can be solved independently, parallelism is possible, which leads to acceleration. A recent development in the area of decomposition techniques is the Expected Separable Overapproximation (ESO) of Richtárik and Takáč [30] and the associated parallel coordinate descent method (PCDM).

Another approach that uses the augmented Lagrangian to tackle the difficulty caused by the linking constraints is the alternating-directions method of multipliers (ADMM) [7, 8]. The ADMM updates the primal variables and Lagrange multipliers using Jacobi iterations, without explicitly optimizing the augmented Lagrangian for any given multiplier. The ADMM is used for stochastic programming problems and referred as progressive hedging [34]. The convergence of ADMM has long been established for two-block problems, however the convergence for multiple blocks has only been addressed recently under specific assumptions [14] or through modifications [12, 44]. We do not present a further discussion of ADMM methods as our focus in this work is on comparing decomposition methods minimizing the augmented Lagrangian for a fixed multiplier.

(Block) coordinate descent methods, early variants of which can be traced back to an 1870 paper of Schwarz [37] and beyond, have recently become very popular due to their low per-iteration cost and good scalability properties. While convergence results were established several decades ago, iteration complexity bounds were not studied until recently [43]. Randomized coordinate and block coordinate descent methods were proposed and analyzed in several settings, such as for smooth convex minimization problems [24, 29, 31], L_1 -regularized problems [?], composite problems [18, 29, 42], nonsmooth convex problems [10], nonconvex problems [19, 25] and problems with separable constraints [22, 23]. Parallel coordinate descent methods were developed and analyzed in [6, 9, 30, 38, 41], primal-dual methods in [39, 41] and inexact methods in [42]. The methods are used in a number of applications, including linear classification [5, 15, 41], compressed sensing [17], truss topology design [28], solving linear systems of equations [16] and group lasso problems [27].

1.1 Augmented Lagrangian

Our work is motivated by the need of solving huge scale instances of constrained convex optimization problems of the form

$$\min_{x^{(1)}, \dots, x^{(n)}} \sum_{i=1}^n g_i(x^{(i)}) \quad (1a)$$

$$\text{subject to } \sum_{i=1}^n A_i x^{(i)} = b \quad (1b)$$

$$x^{(i)} \in X_i, \quad i = 1, 2, \dots, n, \quad (1c)$$

where for $i = 1, 2, \dots, n$ we assume that $X_i \subseteq \mathbf{R}^{N_i}$ are convex and closed sets, $g_i : \mathbf{R}^{N_i} \rightarrow \mathbf{R} \cup \{+\infty\}$ are convex and closed extended real-valued functions and $A_i \in \mathbf{R}^{m \times N_i}$. Problem formulation (1) is also known as an *extended monotropic optimization* problem, using the nomenclature of [2].

While the objective function (1a) and the constraints (1c) are separable in the decision vectors $x^{(1)}, \dots, x^{(n)}$, the linear constraint (1b) links them together, which makes the problem difficult to solve. Moreover, we are interested in the case when n is very large (millions, billions and more), which introduces further computational challenges.

It will be useful to think of the decision vectors $\{x^{(i)}\}$ as ‘blocks’ of a single decision vector $x \in \mathbf{R}^N$, with $N = \sum_i N_i$. This can be achieved as follows. We first partition the $N \times N$ identity matrix I columnwise into n submatrices $U_i \in \mathbf{R}^{N \times N_i}$, $i = 1, 2, \dots, n$, so that $I = [U_1, \dots, U_n]$, and then set $x = \sum_{i=1}^n U_i x^{(i)}$. That is, x is the vector composed by stacking the vectors $x^{(i)}$ on top of each other. It is easy to see that $x^{(i)} = U_i^T x \in \mathbf{R}^{N_i}$. Moreover, if we let $A \stackrel{\text{def}}{=} \sum_{i=1}^n A_i U_i^T \in \mathbf{R}^{m \times N}$, then (1b) can be written compactly as $Ax = b$. Note also that $A_i = AU_i$, for $i = 1, 2, \dots, n$. If we now write $g(x) \stackrel{\text{def}}{=} \sum_i g_i(x^{(i)})$ and $X \stackrel{\text{def}}{=} \{x = \sum_i U_i x^{(i)} | x^{(i)} \in X_i\} \subseteq \mathbf{R}^N$, then problem (1a)–(1c) takes the following form:

$$\begin{aligned} \min_{x \in X} \quad & g(x) \\ \text{subject to} \quad & Ax = b. \end{aligned} \tag{2a} \tag{2b}$$

A typical approach to overcome the issue of nonseparability of the linking constraint (2b) is to drop it and instead consider the *augmented Lagrangian*,

$$F_\pi(x) \stackrel{\text{def}}{=} g(x) + \langle \pi, b - Ax \rangle + \frac{r}{2} \|b - Ax\|^2,$$

where $\pi \in \mathbf{R}^m$ is a Lagrange multiplier vector, $r > 0$ is a penalty parameter and $\|\cdot\|$ denotes the standard Euclidean norm. Now, the Method of Multipliers [3, 13] can be employed to solve problem (1) as described below (Algorithm 1).

Algorithm 1 (Method of Multipliers)

- 1: **Initialization:** $\pi_0 \in \mathbf{R}^m$ and iteration counter $k = 0$
- 2: **while** the stopping condition has not been met **do**
- 3: **Step 1:** Fix the multiplier π_k and solve

$$z_k \leftarrow \min_{x \in X} F_{\pi_k}(x). \tag{3a}$$

- 4: **Step 2:** Update the multiplier

$$\pi_{k+1} \leftarrow \pi_k + r(b - Az_k), \tag{3b}$$

and update the iteration counter $k \leftarrow k + 1$.

- 5: **end while**
-

2 The Problem and Our Contributions

The focus of this paper is on the optimization problem (3a). Hence, we need not be concerned about the dependence of F on π and will henceforth refer to the objective function, dropping the

constant term $\langle \pi, b \rangle$, as $F(x)$. Ignoring the constant term $\langle \pi, b \rangle$, problem (3a) is a *convex composite* optimization problem, i.e., a problem of the form

$$\min_{x \in \mathbf{R}^N} \{F(x) \stackrel{\text{def}}{=} f(x) + \Psi(x)\}, \quad (4)$$

where f is a smooth convex function and Ψ is a separable (possibly nonsmooth) convex function. Indeed, we may set

$$f(x) \stackrel{\text{def}}{=} \frac{r}{2} \|b - Ax\|^2 = \frac{r}{2} \|b - \sum_{i=1}^n A_i x^{(i)}\|^2, \quad (5)$$

and

$$\Psi(x) \stackrel{\text{def}}{=} \begin{cases} g(x) - \langle \pi, Ax \rangle, & x \in X, \\ +\infty, & \text{otherwise.} \end{cases}$$

The main purpose of this work is to draw links between two existing decomposition methods for solving (4), one old and one new, both based on separable approximations to the objective function. In particular, we consider DQAM of Mulvey and Ruszczyński [20, 21, 36] and PCDM of Richtárik and Takáč [30]. Our main contributions (not in order of significance) include:

1. **Two measures of separability.** We show that the parameter “number of neighbours”, used in the analysis of DQAM [36], and the degree of partial separability, used in the analysis of PCDM [30], coincide up to an additive constant in the case of quadratic f .
2. **Two generalizations of DQAM.** We provide a simplified derivation of the diagonal quadratic approximation, which enables us to propose two generalizations of DQAM (Section 3.2) to non-quadratic functions f , based on
 - (i) a finite difference separable approximation of the augmented Lagrangian (Algorithm 3), and
 - (ii) a quadratic approximation with the Hessian matrix replaced by an approximation of its block diagonal (Algorithm 4).

Studying the complexity of these algorithms is outside the scope of this work.

3. **Equivalence of PCDM and DQAM for smooth problems.** We identify a situation in which the second of our generalizations of DQAM (Algorithm 4) coincides with a “fully parallel” variant of PCDM for an appropriate selection of algorithm parameters (see Section 5.4, Theorem 7). This happens for problems with arbitrary smooth f and $\Psi \equiv 0$.
4. **Improved complexity of PCDM under strong convexity.** We derive an improved complexity result for PCDM in the case when F is strongly convex (see Section 6, Theorem 10). The result is much better than that in [30] in situations where the strong convexity constant of F is much larger than the sum of the strong convexity constants of the f and Ψ .
5. **Comparison of the convergence rates of DQAM and PCDM.** We study the newly developed complexity guarantees for (a fully parallel variant of) PCDM and the existing convergence rates for DQAM and show that even though DQAM is specifically designed

to approximate the augmented Lagrangian, PCDM has much better theoretical guarantees (Section 6.3). In particular, if F is strongly convex, both DQAM and PCDM converge linearly; that is, $F(x_{k+1}) \leq qF(x_k)$, where q depends on the method. However, we show that q is much better (i.e., smaller) for PCDM than for DQAM, which then leads to vast speedups in terms of iteration complexity. In particular, we show that the theoretical bound for the number of iterations required to find an ϵ -approximate solution is at least

$$\frac{16(\omega - 1)^3}{\omega} \times \frac{L'}{\bar{L}} \quad (\geq 8 \frac{L'}{\bar{L}} (\omega - 1)^2 \text{ for } \omega \geq 2) \quad (6)$$

times larger for DQAM than for (fully parallel) PCDM. Here, ω is the degree of partial separability¹ of f (defined in Section 4.4), and L' and \bar{L} are the maximum and average of the constants $L_i = r\|A_i^T A_i\|$, $i = 1, 2, \dots, n$, respectively. Note that the speedup factor (6) is larger than 1000 for $\omega = 10$ even in the case when $L' = \bar{L}$. In practice, however, L' will typically be larger than \bar{L} , often much larger.

The form of the speedup factor (6) comes from the fact that DQAM depends on $(\omega - 1)^3$ and L' (while PCDM depends on ω and \bar{L}), which adversely affects its theoretical complexity rate. Let us comment that Mulvey and Ruszczyński [20] remarked that the dependence of DQAM on ω is in practice much better than cubic, although this was not previously established theoretically. We thus answer their conjecture in the affirmative, albeit for a (as we shall see, not so very) different method. To the best of our knowledge, no improved results were available in the literature up to this point.

6. Optimal number of block updates per iteration. PCDM is flexible in that it allows for an *arbitrary* number of block updates per iteration, whereas DQAM needs to update *all blocks*. Moreover, we show that under a simple parallel computing model it is optimal for PCDM to update as many blocks in a single iteration as there are parallel processors (Section 6.4, Theorem 13). As a consequence, the DQAM approach of updating *all* blocks in a single iteration is less than optimal.

7. Computations. We also provide preliminary numerical results that show the practical advantages of PCDM.

We also comment that, as shown in [29], PCDM enjoys complexity guarantees even in the case when F is merely convex, as opposed to it being strongly convex.

3 Diagonal Quadratic Approximation Method

In this section we present the Diagonal Quadratic Approximation Method (DQAM) that was introduced and analysed in a series of papers by Mulvey and Ruszczyński [20, 21], Ruszczyński [36] and Berger, Mulvey and Ruszczyński [1]. As explained in Section 1.1, the augmented Lagrangian is nonseparable because of the cross products $\langle A_i h^{(i)}, A_j h^{(j)} \rangle$ appearing in $f(x + h)$. The DQAM provides a separable approximation of $f(x + h)$ by ignoring these cross terms; this approximation is referred to as the diagonal quadratic approximation (DQA). This makes Step 1 of the method of multipliers ((3a) in Algorithm 1) significantly easier to solve, and amenable to parallel processing.

¹The multiplicative improvement factor (6) is only valid for $\omega \geq 2$ as DQAM was not analyzed in the case $\omega = 1$.

First, notice that we can write

$$\begin{aligned}
f(x+h) &= \frac{r}{2} \|b - A(x+h)\|^2 \\
&= \frac{r}{2} \|b\|^2 - r \langle b, A(x+h) \rangle + \frac{r}{2} (\|Ax\|^2 + 2 \langle Ax, Ah \rangle + \|Ah\|^2) \\
&= f(x) + \langle f'(x), h \rangle + \frac{r}{2} \|Ah\|^2 \\
&= f(x) + \langle f'(x), h \rangle + \frac{r}{2} \sum_{i=1}^n \|A_i h^{(i)}\|^2 + \frac{r}{2} (\|Ah\|^2 - \sum_{i=1}^n \|A_i h^{(i)}\|^2) \\
&= f(x) + \sum_{i=1}^n \langle (f'(x))^{(i)}, h^{(i)} \rangle + \frac{r}{2} \sum_{i=1}^n \|A_i h^{(i)}\|^2 + \frac{r}{2} \sum_{i \neq j} \langle A_i h^{(i)}, A_j h^{(j)} \rangle. \tag{7}
\end{aligned}$$

Now observe that it is only the last term in (7), composed of products $\langle A_i h^{(i)}, A_j h^{(j)} \rangle$ for $i \neq j$, which is not separable. Ignoring these terms, we get a separable approximation of $f(x+h)$ in h ,

$$f(x+h) \approx f^{\text{DQA}}(x+h) \stackrel{\text{def}}{=} f(x) + \langle f'(x), h \rangle + \frac{r}{2} \sum_{i=1}^n \|A_i h^{(i)}\|^2, \tag{8}$$

which in turn leads to a separable approximation of $F(x+h)$ in h :

$$F(x+h) \stackrel{(4)}{=} f(x+h) + \Psi(x+h) \stackrel{(8)}{\approx} f^{\text{DQA}}(x+h) + \Psi(x+h). \tag{9}$$

Mulvey and Ruszczyński [20] propose a slightly less transparent construction of the same approximation. For a fixed x , they approximate $f(y)$ via replacing the cross-products $\langle A_i y^{(i)}, A_j y^{(j)} \rangle$, for $i \neq j$, by

$$\langle A_i y^{(i)}, A_j x^{(j)} \rangle + \langle A_i x^{(i)}, A_j y^{(j)} \rangle - \langle A_i x^{(i)}, A_j x^{(j)} \rangle. \tag{10}$$

Clearly, this is equivalent to what we do above, which can be verified by substituting $y = x+h$ into (10).

3.1 The algorithm

We now present the DQA method (Algorithm 2). The algorithm replaces Step 1 of the Method of Multipliers (Algorithm 1). In what follows, $\theta \in (0, 1)$ is a user defined parameter.

Algorithm 2 (DQAM: Diagonal Quadratic Approximation Method)

1: **for** $k = 0, 1, 2, \dots$ **do**

2: **Step 1a:** Solve for h_k

$$h_k \leftarrow \arg \min_{h \in \mathbf{R}^N} \{f^{\text{DQA}}(x_k + h) + \Psi(x_k + h)\} \tag{11a}$$

3: **Step 1b:** Determine intermediate vector y_k

$$y_k \leftarrow x_k + h_k \tag{11b}$$

4: **Step 1c:** Form the new iterate x_{k+1}

$$x_{k+1} \leftarrow (1 - \theta)x_k + \theta y_k \tag{11c}$$

5: **end for**

Let us now comment on the individual steps of Algorithm 2. Step 1a is easy to execute because the function that is being minimized in (11a) is separable in h , and hence the problem decomposes into n independent *lower-dimensional* problems:

$$h_k^{(i)} = \arg \min_{h^{(i)} \in \mathbf{R}^{N_i}} \left\{ \langle (f'(x_k))^{(i)}, h^{(i)} \rangle + \frac{r}{2} \|A_i h^{(i)}\|^2 + \Psi_i(x_k^{(i)} + h^{(i)}) \right\}, \quad i = 1, 2, \dots, n.$$

Moreover, the problems are *independent*, and hence the updates $h_k^{(1)}, \dots, h_k^{(n)}$ can be computed *in parallel*. In (11b) an intermediate vector y_k is formed, and then in (11c) a convex combination of the current iterate x_k and the intermediate vector y_k is taken to produce the new iterate x_{k+1} . Step (11c) is needed because DQAM uses a local approximation, so if the new point $x_k + h_k$ is far from x_k , the approximation error may be too big and a reduction in the objective function value is not guaranteed. This would lead to serious stability and convergence problems in general, and hence, Step 1c is employed as a correction step for regularizing the method.

3.2 Two generalizations

DQAM was originally designed and analyzed for convex quadratics. Here we propose two generalizations of the method to non-quadratic convex functions f . Our generalizations are based on the following simple result.

Proposition 1. *Let $f(x) = \frac{r}{2} \|b - Ax\|^2$, and let $C_i(x) = U_i^T f''(x) U_i$, where $f''(x)$ is the second derivative (Hessian) of $f(x)$. Then for all $x, h \in \mathbf{R}^N$,*

$$f^{DQA}(x + h) = f(x) + \sum_{i=1}^n \left[f(x + U_i h^{(i)}) - f(x) \right] \quad (12)$$

and

$$f^{DQA}(x + h) = f(x) + \sum_{i=1}^n \left[\langle (f'(x))^{(i)}, h^{(i)} \rangle + \frac{1}{2} \langle C_i(x) h^{(i)}, h^{(i)} \rangle \right]. \quad (13)$$

Proof. First note that

$$\begin{aligned} \sum_{i=1}^n \left[f(x + U_i h^{(i)}) - f(x) \right] &= \sum_{i=1}^n \left[\frac{r}{2} \|b - Ax - A_i h^{(i)}\|^2 - \frac{r}{2} \|b - Ax\|^2 \right] \\ &= \sum_{i=1}^n \left[r \langle Ax - b, A_i h^{(i)} \rangle + \frac{r}{2} \|A_i h^{(i)}\|^2 \right] \\ &= \langle f'(x), h \rangle + \frac{r}{2} \sum_{i=1}^n \|A_i h^{(i)}\|^2, \end{aligned}$$

which, in view of (8), establishes (12). Finally, (13) follows from (8) and the fact that $\frac{r}{2} \|A_i h^{(i)}\|^2 = \frac{1}{2} \langle U_i^T f''(x) U_i h^{(i)}, h^{(i)} \rangle$, which in turn follows from the identities $f''(x) = r A^T A$ and $A_i = A U_i$. \square

Our two generalized methods are obtained by replacing $f^{DQA}(x + h)$ in Step 1 of Algorithm 2 by one of the two approximations (12) and (13). (In the second case, $C_i(x)$ is an arbitrary positive semidefinite matrix and not necessarily $U_i^T f''(x) U_i$), leading to Algorithms 3 and 4, respectively.

Algorithm 3 (Generalization of DQAM: Finite Differences Approximation)

1: **for** $k = 0, 1, 2, \dots$ **do**

2: **Step 1a:** Solve for h_k

$$h_k \leftarrow \arg \min_{h \in \mathbf{R}^N} \left\{ f(x_k) + \sum_{i=1}^n \left[f(x_k + U_i h^{(i)}) - f(x_k) \right] + \Psi(x_k + h) \right\} \quad (14a)$$

3: **Step 1b:** Determine intermediate vector y_k

$$y_k \leftarrow x_k + h_k \quad (14b)$$

4: **Step 1c:** Form the new iterate x_{k+1}

$$x_{k+1} \leftarrow (1 - \theta)x_k + \theta y_k \quad (14c)$$

5: **end for**

Algorithm 3 is based on a finite difference approximation, and is applicable to (possibly) non-smooth functions. Algorithm 4 is based on a separable quadratic approximation. To the best of our knowledge, these algorithms have not been previously proposed, with the exception of the case when f is a convex quadratic when both methods coincide with DQAM.

Algorithm 4 (Generalization of DQAM: Separable Quadratic Approximation)

1: **for** $k = 0, 1, 2, \dots$ **do**

2: **Step 1a:** Solve for h_k

$$h_k \leftarrow \arg \min_{h \in \mathbf{R}^N} \left\{ f(x_k) + \langle f'(x_k), h \rangle + \frac{1}{2} \sum_{i=1}^n \langle C_i(x_k) h^{(i)}, h^{(i)} \rangle + \Psi(x_k + h) \right\} \quad (15a)$$

3: **Step 1b:** Determine intermediate vector y_k

$$y_k \leftarrow x_k + h_k \quad (15b)$$

4: **Step 1c:** Form the new iterate x_{k+1}

$$x_{k+1} \leftarrow (1 - \theta)x_k + \theta y_k \quad (15c)$$

5: **end for**

Algorithm 3 is a derivative-free method that is based upon a finite differences approximation to f . Therefore, this method may be particularly useful when computing derivatives of f is prohibitively expensive, which may be the case for non-quadratic f .

Algorithm 4 replaces the (block of the) Hessian of f , with an approximation to it. For quadratic f , the Hessian is fixed throughout the algorithm, whereas the Hessian changes at each iteration in the non-quadratic case. Therefore, Algorithm 4 may be particularly useful when access to the Hessian is not possible, or when using an approximation to the Hessian is much cheaper than using the true Hessian. We may think of the generalization from DQAM to Algorithm 4, as an

analogue with a Newton method versus a quasi-Newton method. There is an abundance of literature advocating the use of quasi-Newton methods, and so, intuitively, Algorithm 4 may be particularly useful when f is not quadratic.

The convergence analysis for these two methods is an open problem, which we do not consider in this paper. Instead, we propose that DQAM be replaced by PCDM, described in the next section.

4 Parallel Coordinate Descent Method

As discussed in the introduction, we propose that instead of implementing Step 1 of the Method of Multipliers (Algorithm 1) using DQAM, a parallel coordinate descent method (PCDM) be used. This section is devoted to describing the method, developed by Richtárik and Takáč [30].

4.1 Block samplings

As we shall see, unlike DQAM where all blocks are updated at each iteration, PCDM allows for an (almost) arbitrary random subset of blocks to be updated at each iteration. The purpose of this section is to formalize this.

In particular, at iteration k only blocks $i \in S_k \subseteq \{1, 2, \dots, n\}$ are updated, where $\{S_k\}$, $k \geq 0$, are iid random sets having the following two properties:

$$\mathbf{P}(i \in S_k) = \mathbf{P}(j \in S_k) \quad \text{for all } i, j \in \{1, 2, \dots, n\}, \quad (16)$$

$$\mathbf{P}(i \in S_k) > 0 \quad \text{for all } i \in \{1, 2, \dots, n\}. \quad (17)$$

It is clear that, necessarily, $\mathbf{P}(i \in S_k) = \frac{\mathbf{E}[\|S_k\|]}{n}$. Following [30], for simplicity we refer to an arbitrary random set-valued mapping with values in the power set $2^{\{1, 2, \dots, n\}}$ by the name *block sampling*, or simply *sampling*. A sampling S_k is called *uniform* if it satisfies (16) and *proper* if it satisfies (17).

In [30], PCDM was analyzed for all proper uniform samplings. However, better complexity results were obtained for so called *doubly uniform* samplings, which belong to the family of uniform samplings. For brevity purposes, in this paper we concentrate on a subclass of doubly uniform samplings called τ -nice samplings, which we now define.

Definition 2 (τ -nice sampling). *Let τ be an integer between 1 and n . A sampling \hat{S} is called τ -nice if for all $S \subseteq \{1, 2, \dots, n\}$,*

$$\mathbf{P}(\hat{S} = S) = \begin{cases} 0, & |S| \neq \tau, \\ \frac{1}{\binom{n}{\tau}}, & \text{otherwise.} \end{cases}$$

(Here, $1/\binom{n}{\tau}$ denotes ‘ n choose τ ’.) A natural candidate for τ is the number of available processors/threads as then updates to the τ blocks of x_k can be computed in parallel. As we shall later see, this is also the optimal choice from the complexity point of view (Theorem 13).

4.2 Expected Separable Overapproximation (ESO)

Fix positive scalars w_1, \dots, w_n . (We write $w = (w_1, \dots, w_n)$.) We define a separable norm by

$$\|x\|_w^2 \stackrel{\text{def}}{=} \sum_{i=1}^n w_i \|x^{(i)}\|_{(i)}^2, \quad x \in \mathbf{R}^N, \quad (18)$$

where for each $i = 1, 2, \dots, n$ we fix a positive definite matrix $B_i \in \mathbf{R}^{N_i \times N_i}$ and set

$$\|t\|_{(i)} \stackrel{\text{def}}{=} \langle B_i t, t \rangle^{1/2}, \quad t \in \mathbf{R}^{N_i}. \quad (19)$$

We can now define the concept of expected separable overapproximation.

Definition 3 (Expected Separable Overapproximation (ESO) [30]). *Let β and w_1, \dots, w_n be positive constants and \hat{S} be a proper uniform sampling. We say that $f : \mathbf{R}^N \rightarrow \mathbf{R}$ admits a (β, w) -ESO with respect to \hat{S} (and, for simplicity, we write $(f, \hat{S}) \sim \text{ESO}(\beta, w)$) if for all $x, h \in \mathbf{R}^N$,*

$$\mathbf{E} \left[f \left(x + \sum_{i \in \hat{S}} U_i h^{(i)} \right) \right] \leq f(x) + \frac{\mathbf{E}[\|\hat{S}\|]}{n} \left(\langle f'(x), h \rangle + \frac{\beta}{2} \|h\|_w^2 \right). \quad (20)$$

In Section 4.3 we describe how the ESO is used to design a parallel coordinate descent method for solving problem (4). The issue of how the parameters w and β giving rise to an ESO can be determined/computed will be discussed in Section 4.5.

4.3 The algorithm

In PCDM we compute the approximation

$$\mathbf{E} \left[F(x + \sum_{i \in \hat{S}} U_i h^{(i)}) \right] = \mathbf{E} \left[f(x + \sum_{i \in \hat{S}} U_i h^{(i)}) + \Psi(x + \sum_{i \in \hat{S}} U_i h^{(i)}) \right]. \quad (21)$$

It can be verified (see [30, Section 3]) that due to separability of Ψ the following identity holds:

$$\mathbf{E} \left[\Psi(x + \sum_{i \in \hat{S}} U_i h^{(i)}) \right] = \left(1 - \frac{\mathbf{E}[\|\hat{S}\|]}{n} \right) \Psi(x) + \frac{\mathbf{E}[\|\hat{S}\|]}{n} \Psi(x + h). \quad (22)$$

Substituting (22) and (20) into (21), we obtain

$$\mathbf{E} \left[F(x + \sum_{i \in \hat{S}} U_i h^{(i)}) \right] \leq F^{\text{ESO}}(x + h) \stackrel{\text{def}}{=} \left(1 - \frac{\mathbf{E}[\|\hat{S}\|]}{n} \right) F(x) + \frac{\mathbf{E}[\|\hat{S}\|]}{n} H_{\beta, w}(x + h), \quad (23)$$

where

$$H_{\beta, w}(x + h) \stackrel{\text{def}}{=} f(x) + \langle f'(x), h \rangle + \frac{\beta}{2} \|h\|_w^2 + \Psi(x + h), \quad (24)$$

which is separable in h :

$$H_{\beta, w}(x + h) \stackrel{(18)+(19)}{=} f(x) + \sum_{i=1}^n \left\{ \langle (f'(x))^{(i)}, h^{(i)} \rangle + \frac{\beta w_i}{2} \langle B_i h^{(i)}, h^{(i)} \rangle + \Psi_i(x^{(i)} + h^{(i)}) \right\}. \quad (25)$$

We are now ready to present the parallel coordinate descent method (Algorithm 5).

Algorithm 5 (PCDM: Parallel Coordinate Descent Method)

- 1: **Initialization:** $x_0 \in \mathbf{R}^N$, ESO parameters (β, w)
- 2: **for** $k = 0, 1, 2, \dots$ **do**
- 3: **Step 1a:** Solve

$$h_k \leftarrow \arg \min_{h \in \mathbf{R}^N} F^{\text{ESO}}(x_k + h) \quad (26a)$$

- 4: **Step 1b:** Update x_k

$$x_{k+1} \leftarrow x_k + \sum_{i \in S_k} U_i h_k^{(i)} \quad (26b)$$

- 5: **end for**
-

Given an iterate x_k , in (26a) we compute

$$h_k = h(x_k) \stackrel{\text{def}}{=} \arg \min_{h \in \mathbf{R}^N} F^{\text{ESO}}(x_k + h) \stackrel{(23)}{=} \arg \min_{h \in \mathbf{R}^N} H_{\beta, w}(x_k + h). \quad (27)$$

Further, note that (26b) is equivalent to writing $x_{k+1}^{(i)} = x_k^{(i)} + h_k^{(i)}$ for all $i \in S_k$ and $x_{k+1}^{(i)} = x_k^{(i)}$ for all $i \notin S_k$. That is, only blocks belonging to the random set S_k are updated. This means that in (26a) we need not compute all blocks of h_k . In view of (25) and (27), this is possible, and hence (26a) can be replaced by

$$h_k^{(i)} \leftarrow \arg \min_{h^{(i)} \in \mathbf{R}^{N_i}} \left\{ \langle (f'(x_k))^{(i)}, h^{(i)} \rangle + \frac{\beta w_i}{2} \langle B_i h^{(i)}, h^{(i)} \rangle + \Psi_i(x_k^{(i)} + h^{(i)}) \right\}, \quad i \in S_k. \quad (28)$$

4.4 Measure of separability

In this section we provide a link between the measures of separability of f utilized in the analysis of DQAM [20] and PCDM [30]. In the first case, the quantity is defined specifically for a quadratic objective; in the second case the definition is general. As we shall see, both quantities coincide in the quadratic case. As the complexity of the two methods depends on these quantities, our observation allows us to compare the convergence rates (see Section 6). Both measures of separability are to be understood with respect to the fixed block structure introduced before.

Now we define the degree of partial separability introduced by Richtárik and Takáč [30] for a smooth convex function.

Definition 4 (Partial separability). *A smooth convex function $f : \mathbf{R}^N \rightarrow \mathbf{R}$ is partially separable of degree ω if there exists a collection \mathcal{J} of subsets of $\{1, 2, \dots, n\}$ such that*

$$f(x) = \sum_{J \in \mathcal{J}} f_J(x) \quad \text{and} \quad \max_{J \in \mathcal{J}} |J| \leq \omega, \quad (29)$$

where for each J , f_J is a smooth convex function that depends on $x^{(i)}$ for $i \in J$ only.

The following result establishes the relationship between ω and the *Ruszczyński measure of separability*, ω_R , defined in [36].² This will be important in Section 6 for the complexity comparison.

Theorem 5. *For convex quadratic function f given by (5) we have $\omega = \omega_R + 1$, where ω_R is the Ruszczyński measure of separability defined in [36].*

Proof. First, we can write

$$f(x) = \frac{r}{2} \sum_{j=1}^m \left(b_j - \sum_{i=1}^n A_{ji} x^{(i)} \right)^2, \quad (30)$$

where b_j is the j -th entry of b . Note that all summands in the decomposition are convex and smooth. Moreover, summand j depends on $x^{(i)}$ if and only if $A_{ji} \neq 0$. If we now let

$$\omega_j = |\{i : A_{ji} \neq 0\}|, \quad j = 1, 2, \dots, m, \quad (31)$$

²Note that, in [36], the *Ruszczyński measure of separability* is denoted by N , but we call it ω_R here to avoid a notation clash. See the notation dictionary in Appendix A for further details. Furthermore, ω_R is defined for a strongly convex quadratic function only.

then we conclude that f is partially separable of degree

$$\omega = \max_{j \in \{1, 2, \dots, m\}} \omega_j. \quad (32)$$

In words, ω is the maximum number of blocks linked by any single constraint. Moreover, “ ω_R is the maximum number of blocks linked by any single constraint, decremented by one” [36]. Thus establishing the result. (Further details on ω_R can be found in [36] and in Appendix B.) \square

4.5 ESO for partially separable smooth convex functions

In order for PCDM to be implementable, one needs first to compute the parameters w_1, \dots, w_n (defining the norm $\|\cdot\|_w$) and $\beta > 0$ for which $(f, \hat{S}) \sim \text{ESO}(\beta, w)$, i.e., for which (20) holds. Clearly, the parameters β and w depend on f and \hat{S} .

In what follows we will assume that the gradient of f is block Lipschitz. That is, there exist positive constants L_1, \dots, L_n such that for all $x \in \mathbf{R}^N$, $i \in \{1, 2, \dots, n\}$ and $h^{(i)} \in \mathbf{R}^{N_i}$,

$$\|(f'(x + U_i t))^{(i)} - (f'(x))^{(i)}\|_{(i)}^* \leq L_i \|t\|_{(i)}, \quad (33)$$

where $\|s\|_{(i)}^* \stackrel{\text{def}}{=} \max\{\langle s, x \rangle : \|x\|_w = 1\} = \langle B_i^{-1} s, s \rangle^{1/2}$ is the conjugate norm to $\|\cdot\|_w$.

Theorem 6 (Theorem 14 in [30]). *Assume f is convex, partially separable of degree ω , and has block Lipschitz gradient with constants $L_1, L_2, \dots, L_n > 0$. Further, assume that \hat{S} is a τ -nice sampling, where $\tau \in \{1, 2, \dots, n\}$. Then $(f, \hat{S}) \sim \text{ESO}(\beta, w)$, where*

$$\beta = 1 + \frac{(\omega - 1)(\tau - 1)}{\max\{1, n - 1\}}, \quad w_i = L_i, \quad i = 1, 2, \dots, n. \quad (34)$$

In Section 6 we study the complexity of PCDM in the case covered by the above theorem (and under a further strong convexity assumption).

Consider now the special case of convex quadratic f given by (5). If the matrices $A_i^T A_i$, $i = 1, 2, \dots, n$, are all positive definite, we can choose $B_i = r A_i^T A_i$, $i = 1, 2, \dots, n$, in which case we will have $L_i = 1$ for all i . Otherwise we can choose B_i to be the $N_i \times N_i$ identity matrix, and then

$$L_i = r \|A_i^T A_i\| \stackrel{\text{def}}{=} r \max_{\|h^{(i)}\| \leq 1} \|A_i^T A_i h^{(i)}\|, \quad (35)$$

where both norms in the definition are the standard Euclidean norms in \mathbf{R}^{N_i} .

4.6 Fully parallel coordinate descent method

PCDM used with an n -nice sampling \hat{S} resembles DQAM in two ways: i) it updates *all* blocks during each iteration, ii) it is not randomized. Indeed, $h = \sum_{i=1}^n U_i h^{(i)} = \sum_{i \in \hat{S}} U_i h^{(i)}$, and hence

$$F(x + h) = \mathbf{E}[F(x + h)] \stackrel{(23)}{\leq} F^{\text{ESO}}(x + h) \stackrel{(23)+(24)}{=} f(x) + \langle f'(x), h \rangle + \frac{\beta}{2} \|h\|_w^2 + \Psi(x + h).$$

In particular, in the setting of Theorem 6 we have $\beta = \omega$ and $w = L = (L_1, \dots, L_n)$, and Algorithm 5 specializes to fully parallel PCDM.

5 Links Between DQAM and PCDM

In this section we discuss and compare DQAM and PCDM. We highlight some of the main differences between the two methods, and describe a special case where the methods coincide.

5.1 Fully parallel vs partially parallel updating

One of the main differences between DQAM and PCDM is the number of blocks that must be updated at each iteration. At each iteration of DQAM, *all* n blocks must be updated. This highlights the fact that DQAM uses a *fully parallel* update scheme. On the other hand, PCDM is more flexible as it is able to update τ blocks at each iteration where $1 \leq \tau \leq n$. This is beneficial because in practice there are usually fewer processors than the number of blocks. So, PCDM can act as a *serial* method if $\tau = 1$, a *fully parallel* method if $\tau = n$, or it can be optimized to the number of processors p (so $\tau = p$). The advantages of updating $\tau = p$ blocks at each iteration of PCDM is established theoretically in Section 6.

Because DQAM updates all n blocks at each iteration, it is a Jacobi type method, whereas PCDM can be interpreted as a Jacobi type method when $\tau = n$, a Gauss-Seidel type method when $\tau = 1$, or a hybrid Jacobi-Gauss-Seidel method when $1 < \tau < n$.

5.2 Flexibility of PCDM

PCDM can be applied to a general convex composite function. Specifically, f is only assumed to be smooth and convex. Further, the algorithm is guaranteed to converge when applied to a general smooth convex function, and can be equipped with iteration complexity bounds (see [30]). On the other hand, the convergence results for DQAM have been only derived under the assumption that f is quadratic and strongly convex; there are no convergence guarantees for a function f with any other structure. Complexity estimates for both methods are discussed in detail in Section 6.

Notice that DQAM has been tailored specifically for an augmented Lagrangian objective function so it is reasonable that the function f is assumed to be quadratic and strongly convex in this context. However, this assumption restricts the range of problems that can be solved using DQA, while PCDM can be applied to a much wider class of problems.

5.3 Approximation type and algorithm philosophy

In DQAM, a local two-sided approximation to the cross products is employed. The error associated with the approximation is of the order $o(\|h\|_2^2)$, which explains that, if the update h_k is too large, then the model loses accuracy. This justifies the need for a correction step (11c) so as to ensure that x_{k+1} is not too far from x_k . This ensures a reduction in the objective value and ultimately, algorithm convergence. The need for a correction scheme within DQAM is also apparent from the finite differences formulation presented in Algorithm 3. Consider the summation in (14a), and for simplicity assume that $\Psi \equiv 0$. Then the block update $h_k^{(i)}$ is that which minimizes the function value difference in the i -th block coordinate direction, independently of all the other blocks $j \neq i$. Clearly, this will *not* guarantee that $F(x_k + h_k) \leq F(x_k)$ because the function F is not block separable. A simple 2D quadratic example showing that this approach is doomed to fail was described in [41].

In contrast to the DQAM scheme, PCDM employs a one-sided *global* expected separable overapproximation of the augmented Lagrangian (4), which guarantees to produce a new random iterate

x_{k+1} that, on average, decreases the objective function. That is, x_{k+1} satisfies $\mathbf{E}[F(x_{k+1}) \mid x_k] \leq F(x_k)$. It turns out that this is sufficient to obtain a high probability complexity result and therefore there is no need for a correction step in PCDM. In fact, as we shall see in Section 5.4, a “correction step” is already embedded in the approximation in the form of the ESO parameter β .

Note that, besides DQAM, there are many other algorithms that follow a “step-then-correct” strategy. One example are trust region methods, where a solution to some subproblem is found, the quality of the solution is measured, and then the size of the trust region is adjusted to reflect the quality. A second example is the conditional gradient algorithm, which builds a linear approximation to the objective function, finds the minimizer of the linearized problem (the “step”) and then “corrects” by taking a convex combination of the previous point and the step to reduce the objective value. A correction step is implicitly built-in for PCDM, in the choice of the constant β .

5.4 A special case in which the methods coincide

So far we have highlighted some of the differences between DQAM and PCDM. However, in this section we present a special case where the two methods coincide.

Theorem 7. *Assume f is partially separable of degree ω , and has block Lipschitz gradient with constants $L_1, L_2, \dots, L_n > 0$. Further, assume $\Psi \equiv 0$. Then Algorithm 4 (generalization of DQAM) coincides with fully parallel PCDM under the following choice of parameters:*

$$C_i(x_k) \equiv L_i B_i \quad (i = 1, 2, \dots, n), \quad \theta = \frac{1}{\omega}. \quad (36)$$

Proof. In Algorithm 4 we have $x_{k+1} = (1 - \theta)x_k + \theta(x_k + h_k)$, where

$$h_k = \arg \min_{h \in \mathbf{R}^N} \left\{ \langle f'(x_k), h \rangle + \frac{1}{2} \sum_{i=1}^n \langle C_i(x_k) h^{(i)}, h^{(i)} \rangle \right\}. \quad (37)$$

Due to separability of the objective function in (37) and the choice of parameters (36), we see that $h_k^{(i)} = -\frac{1}{L_i} B_i^{-1} (f'(x_k))^{(i)}$, $i = 1, 2, \dots, n$, and hence

$$x_{k+1}^{(i)} = (1 - \theta)x_k^{(i)} + \theta(x_k^{(i)} + h_k^{(i)}) = x_k^{(i)} - \frac{1}{\omega L_i} B_i^{-1} (f'(x_k))^{(i)}. \quad (38)$$

In fully parallel PCDM we have $x_{k+1} = x_k + h_k$, where

$$h_k = \arg \min_{h \in \mathbf{R}^N} \left\{ \langle f'(x_k), h \rangle + \frac{\omega}{2} \sum_{i=1}^n \langle L_i B_i h^{(i)}, h^{(i)} \rangle \right\}. \quad (39)$$

Using separability of the objective function in (39), we again obtain the same formula (38) for x_{k+1} , establishing the equivalence of the two methods. \square

A few remarks:

- In the context of the original problem (1), the case covered by the above theorem corresponds to a feasibility problem ($\Psi \equiv 0$ means that $g \equiv 0$).
- DQAM was analyzed in [36] only for the parameter θ in the interval $(0, \frac{1}{2(\omega-1)})$. For $\omega > 1$ this leads to *smaller steps* than the PCDM default choice $\theta = \frac{1}{\omega}$, which, intuitively, is expected to result in slower convergence of DQAM.

6 Complexity of DQAM and PCDM under Strong Convexity

In this section we study and compare the convergence rates of DQAM and PCDM under the assumption of strong convexity of the objective function. We limit ourselves to this case as complexity estimates for DQAM are not available otherwise. Both DQAM and PCDM benefit from linear convergence, but the rate is much better for PCDM than for DQA.

Strong convexity. We assume that F is strongly convex with respect to the norm $\|\cdot\|_w$ for some vector of positive weights $w = (w_1, \dots, w_n)$ specified in the results, with (strong) convexity parameter $\mu_F > 0$. A function $\phi : \mathbf{R}^N \rightarrow \mathbf{R} \cup \{+\infty\}$ is strongly convex with respect to the norm $\|\cdot\|_w$ with convexity parameter $\mu_\phi = \mu_\phi(w) \geq 0$ if for all $x, y \in \text{dom } \phi$,

$$\phi(y) \geq \phi(x) + \langle \phi'(x), y - x \rangle + \frac{\mu_\phi}{2} \|y - x\|_w^2, \quad (40)$$

where $\phi'(x)$ is any subgradient of ϕ at x . The case with $\mu_\phi(w) = 0$ reduces to convexity. It will be useful to note that for any $t > 0$,

$$\mu_\phi(tw) = \mu_\phi(w)/t. \quad (41)$$

Strong convexity of F may come from f or Ψ or both and we will write μ_f (resp. μ_Ψ) for the strong convexity parameter of f (resp. Ψ). It is easy to see that

$$\mu_F \geq \mu_f + \mu_\Psi. \quad (42)$$

Note that the strong convexity constant of F can be *arbitrarily larger* than the sum of the strong convexity constants of the functions f and Ψ . Indeed, consider the following simple 2D example ($N = n = 2$): $f(x) = \frac{\mu}{2}(x^{(1)})^2$, $\Psi(x) = \frac{\mu}{2}(x^{(2)})^2$, where $\mu > 0$. Let $\|x\|_w$ be the standard Euclidean norm (i.e., $B_i = 1$ and $w_i = 1$ for $i = 1, 2$). Clearly, neither f nor Ψ is strongly convex ($\mu_f = \mu_\Psi = 0$). However, F is strongly convex with constant $\mu_F = \mu$.

In the rest of the section we will repeatedly use the following simple result.

Lemma 8. Let $\xi_0 > \epsilon > 0$ and $\gamma \in (0, 1)$. If $k \geq \frac{1}{\gamma} \log \left(\frac{\xi_0}{\epsilon} \right)$, then $(1 - \gamma)^k \xi_0 \leq \epsilon$.

Proof. $(1 - \gamma)^k \xi_0 = (1 - \frac{1}{1/\gamma})^{(1/\gamma)(\gamma k)} \xi_0 \leq e^{-\gamma k} \xi_0 \leq e^{-\log(\xi_0/\epsilon)} \xi_0 = \epsilon$. \square

6.1 PCDM

We now derive a new improved complexity result for PCDM. In [30, Theorem 20] the authors prove an iteration complexity bound based on the assumption that $\mu_f + \mu_\Psi > 0$. Here we obtain a new and tighter complexity result under the weaker assumption $\mu_F > 0$. As discussed above, μ_F can be substantially bigger than $\mu_f + \mu_\Psi$, which implies that our complexity bound can be much better.

The following auxiliary result is an improvement on Lemma 17(ii) in [30] and will be used in the proof of our main complexity result.

Lemma 9. If $\mu_F(w) > 0$ and $\beta \geq \mu_f(w)$, then for all $x \in \text{dom } F$

$$H_{\beta, w}(x + h(x)) - F^* \leq \frac{\beta - \mu_f(w)}{\mu_F(w) + \beta - \mu_f(w)} (F(x) - F^*). \quad (43)$$

Proof. Let $\mu_F = \mu_F(w)$ and $\mu_f = \mu_f(w)$. By Lemma 16 in [30], we have

$$H_{\beta,w}(x + h(x)) \leq \min_{y \in \mathbf{R}^N} \left\{ F(y) + \frac{\beta - \mu_f}{2} \|y - x\|_w^2 \right\}. \quad (44)$$

Using this, we can further write

$$\begin{aligned} H_{\beta,w}(x + h(x)) &\stackrel{(44)}{\leq} \min_{y = \lambda x^* + (1-\lambda)x, \lambda \in [0,1]} \left\{ F(y) + \frac{\beta - \mu_f}{2} \|y - x\|_w^2 \right\} \\ &= \min_{\lambda \in [0,1]} \left\{ F(\lambda x^* + (1-\lambda)x) + \frac{(\beta - \mu_f)\lambda^2}{2} \|x - x^*\|_w^2 \right\} \\ &\leq \min_{\lambda \in [0,1]} \left\{ \lambda F^* + (1-\lambda)F(x) - \frac{\mu_F \lambda(1-\lambda) - (\beta - \mu_f)\lambda^2}{2} \|x - x^*\|_w^2 \right\}, \end{aligned} \quad (45)$$

where in the last step we have used strong convexity of F . Notice that $\lambda^* \stackrel{\text{def}}{=} \mu_F / (\mu_F + \beta - \mu_f) \in (0, 1]$ and that $\mu_F(1 - \lambda^*) - (\beta - \mu_f)\lambda^* = 0$. It now only remains to substitute λ^* into (45) and subtract F^* from the resulting inequality. \square

We now present our main complexity result. It gives a bound on the number of iterations required by PCDM (Algorithm 5) to obtain an ϵ solution with high probability. The result is generic in the sense that it applies to any smooth convex function and proper uniform sampling as long as the parameters β and w giving rise to an ESO are known.

Theorem 10. *Assume that $F = f + \Psi$ is strongly convex with respect to the norm $\|\cdot\|_w$ ($\mu_F(w) > 0$). Let \hat{S} be a proper uniform sampling satisfying $(f, \hat{S}) \sim \text{ESO}(\beta, w)$. Choose an initial point $x_0 \in \mathbf{R}^N$, target confidence level $\rho \in (0, 1)$, target accuracy level $0 < \epsilon < F(x_0) - F^*$ and iteration counter*

$$K \geq \frac{n}{\mathbf{E}[\|\hat{S}\|]} \frac{\beta + \mu_F(w) - \mu_f(w)}{\mu_F(w)} \log \left(\frac{F(x_0) - F^*}{\epsilon \rho} \right). \quad (46)$$

If $\{x_k\}$, $k \geq 0$, are the random points generated by PCDM (Algorithm 5) as applied to problem (4), then $\mathbf{P}(F(x_K) - F^ \leq \epsilon) \geq 1 - \rho$.*

Proof. Let $\alpha = \frac{\mathbf{E}[\|\hat{S}\|]}{n}$ and $\xi_k = F(x_k) - F^*$. Then for all $k \geq 0$,

$$\mathbf{E}[\xi_{k+1} \mid x_k] \stackrel{(23)}{\leq} (1 - \alpha)\xi_k + \alpha(H_{\beta,w}(x_k + h(x_k)) - F^*) \stackrel{(\text{Lemma 9})}{\leq} (1 - \gamma)\xi_k. \quad (47)$$

where $\gamma = \alpha\mu_F(w)/(\mu_F(w) + \beta - \mu_f(w))$. Note that Lemma 9 is applicable as the assumption $\beta \geq \mu_f(w)$ is satisfied due to the fact that $(f, \hat{S}) \sim \text{ESO}(\beta, w)$ (see [30, Section 4]). Further, note that $\gamma > 0$ since $\alpha > 0$ and $\mu_F(w) > 0$. Moreover, $\gamma \leq 1$ since $\alpha \leq 1$ and $\beta \geq \mu_f(w)$. By taking expectation in x_k through (47), we obtain $\mathbf{E}[\xi_k] \leq (1 - \gamma)^k \xi_0$. Applying Markov inequality, Lemma 8 and (46), we obtain $\mathbf{P}(\xi_K > \epsilon) \leq \mathbf{E}[\xi_K]/\epsilon \leq (1 - \gamma)^K \xi_0/\epsilon \leq \rho$, establishing the result. \square

In order to compare the complexity of PCDM with that of DQAM, which is a fully parallel method, we now derive a specialized complexity result for fully parallel PCDM. The method is no longer stochastic in this situation, i.e., the sequence of vectors $\{x_k\}$, $k \geq 0$, is deterministic. Hence, we give a standard complexity result as opposed to a high probability one. Finally, we make use of the fact that for partially separable functions f , the parameters β and w are known.

Theorem 11. Assume $f : \mathbf{R}^N \rightarrow \mathbf{R}$ is partially separable of degree ω , and has block Lipschitz gradient with constants $L_1, L_2, \dots, L_n > 0$. Further assume that $F = f + \Psi$ is strongly convex with $\mu_F(L) > 0$, where $L = (L_1, \dots, L_n)$. Finally, let $\{x_k\}_{k \geq 0}$ be the sequence generated by fully parallel PCDM. Then for all $k \geq 0$,

$$F(x_{k+1}) - F^* \leq q^{PCDM}(F(x_k) - F^*), \quad (48)$$

where

$$q^{PCDM} = 1 - \frac{\mu_F(L)}{\omega + \mu_F(L) - \mu_f(L)}. \quad (49)$$

Moreover, if we let $\epsilon < F(x_0) - F^*$ and

$$k \geq \frac{1}{1 - q^{PCDM}} \log \left(\frac{F(x_0) - F^*}{\epsilon} \right), \quad (50)$$

then $F(x_k) - F^* \leq \epsilon$.

Proof. Let \hat{S} be the fully parallel sampling, i.e., the n -nice sampling. Applying Theorem 6, we see that $(f, \hat{S}) \sim ESO(\beta, w)$, with $\beta = \omega$ and $w = L$. Following the first part of the proof of Theorem 10, we have $\alpha = 1$ and $\xi_{k+1} \leq (1 - \gamma)\xi_k$, where $\gamma = \mu_F(L)/(\mu_F(L) + \omega - \mu_f(L))$, establishing (48). The second statement follows directly by applying Lemma 8. \square

6.2 DQAM

We now present a complexity result for DQAM, established in [36].

Theorem 12 (Theorem 2 in [36]). Let $f(x) = \frac{r}{2} \|b - Ax\|^2$ be partially separable of degree $\omega > 1$. Assume that $F (= f + \Psi)$ is strongly convex with $\mu_F(e) > 0$, where $e \in \mathbf{R}^n$ is the vector of all ones. Further assume that the sets X_i , $i = 1, \dots, n$, are bounded. Let $\{x_k\}$, $k \geq 0$, be the sequence generated by DQAM (Algorithm 2) with $\theta = \frac{1}{2(\omega-1)}$. Then for all $k \geq 0$,

$$F(x_{k+1}) - F^* \leq q^{DQAM}(F(x_k) - F^*),$$

where

$$q^{DQAM} = 1 - \frac{\mu_F(e)}{16L'(\omega - 1)^3 + 4(\omega - 1)\mu_F(e)}, \quad (51)$$

and $L' \stackrel{\text{def}}{=} \max_{1 \leq i \leq n} r \|A_i\|^2$. Moreover, if we let $\epsilon < F(x_0) - F^*$ and

$$k \geq \frac{1}{1 - q^{DQAM}} \log \left(\frac{F(x_0) - F^*}{\epsilon} \right), \quad (52)$$

then $F(x_k) - F^* \leq \epsilon$.

Ruszczynski analyzed DQAM for a range of parameters θ : $\theta \in (0, 1/(\omega - 1))$ [36, Theorem 1; $\mu = 0$]. However, the choice $\theta = 1/(2(\omega - 1))$ is optimal [36, Eq (5.11)], and the above theorem presents Ruszczynski's result for this optimal choice of the stepsize parameter. A table translating the notation used in this paper and [36] is included in Appendix B.

6.3 Comparison of the Linear Rates of DQAM and PCDM

We now compare the convergence rates q^{DQAM} and q^{PCDM} defined in (51) and (49), respectively, and the resulting iteration complexity guarantees. We will argue that q^{PCDM} can be much better (i.e., smaller) than q^{DQAM} , leading to vastly improved iteration complexity bounds. However, as we shall see, in practice the fully parallel PCDM method and DQAM behave similarly, with PCDM being about twice as fast as DQAM.

Before we start with the comparison, recall from (35) that the gradient of $f(x) = \frac{r}{2}\|b - Ax\|^2$ (i.e., f covered by Theorem 12) is block Lipschitz with constants $L_i = r\|A_i^T A_i\|$, $i = 1, 2, \dots, n$. Hence, $L' = \max_i L_i$, which draws a link between the quantities L_i , $i = 1, 2, \dots, n$, appearing in Theorem 11 and L' appearing in Theorem 12.

- **Identical Lipschitz constants.** Assume now that $L_i = L'$ for all $i = 1, 2, \dots, n$ and let $L = (L_1, \dots, L_n)$, as in Theorem 11. Using (41) we observe that

$$\mu_\phi(L) = \mu_\phi(L'e) = \mu_\phi(e)/L', \quad (53)$$

whence

$$q^{\text{PCDM}} \stackrel{(49)+(53)}{=} 1 - \frac{\mu_F(e)}{L'\omega + \mu_F(e) - \mu_f(e)}. \quad (54)$$

We can now directly compare q^{PCDM} and q^{DQAM} by comparing (54) and (51). Clearly³,

$$16L'(\omega - 1)^3 \geq L'\omega \quad \text{and} \quad 4(\omega - 1)\mu_F(e) \geq \mu_F(e) - \mu_f(e), \quad (55)$$

and hence $q^{\text{PCDM}} \leq q^{\text{DQAM}}$. However, both inequalities in (55) can be very loose, which means that q^{PCDM} can be much better than q^{DQAM} . For instance, in the case when $\mu_F(e) = \mu_f(e)$, we have

$$\frac{1 - q^{\text{PCDM}}}{1 - q^{\text{DQAM}}} = \frac{16L'(\omega - 1)^3 + 4(\omega - 1)\mu_F(e)}{L'\omega} \geq \frac{16(\omega - 1)^3}{\omega}. \quad (56)$$

In view of (50) and (52), this means that the number of DQAM iterations needed to obtain an ϵ -solution is larger than that for PCDM *by at least the multiplicative factor* $16(\omega - 1)^3/\omega$. For instance, the theoretical iteration complexity of DQAM is more than 1000 times worse than that of PCDM for $\omega = 10$.

- **Varying Lipschitz Constants.** If the constants L_1, \dots, L_n are not all equal, it is somewhat difficult to compare the complexity rates as we cannot directly compare the strong convexity constants $\mu_\phi(L)$ and $\mu_\phi(e)$ (for $\phi = F$ and $\phi = f$). What we can do, however, is to at least make sure that the “scaling” is identical in both. Here is what we mean by that. Recall that $\mu_\phi(w)$ is the strong convexity constant of ϕ wrt a *weighted norm* $\|x\|_w$ defined by (18). As we have remarked in (41), if we scale the weights by a positive factor $t > 0$, the corresponding strong convexity constant scales by $1/t$. Hence, $\mu_\phi(L)$ and $\mu_\phi(e)$ cannot be considered comparable unless $\sum_i L_i = \sum_i e_i = n$. Of course, even if this was the case, it is possible that the strong convexity constants might be very different. However, in this case there is at least no reason to suspect a-priori that one might be larger than the other, and hence they are comparable in that sense.

³This holds as long as $\omega > 1$, which is the case covered by Theorem 12 and hence assumed here.

If we let $\bar{L} = \frac{1}{n} \sum_i L_i$ and $w_i = L_i/\bar{L}$ for $i = 1, 2, \dots, n$, then $\sum_i w_i = n$, and hence, as explained above, $\mu_\phi(w) \approx \mu_\phi(e)$. Furthermore, since $w = L/\bar{L}$, and using (41), we have $\mu_\phi(L) = \mu_\phi(\bar{L}w) = \mu_\phi(w)/\bar{L} \approx \mu_\phi(e)/\bar{L}$. This is an analogue of (53) and we can therefore now continue our comparison in the same way as we did for the case with identical Lipschitz constants. In particular, if $\mu_F(e) = \mu_f(e)$ (for simplicity), then as above we can argue that

$$\frac{1 - q^{\text{PCDM}}}{1 - q^{\text{DQAM}}} = \frac{16L'(\omega - 1)^3 + 4(\omega - 1)\mu_F(e)}{\bar{L}\omega} \geq \frac{16(\omega - 1)^3 L'}{\omega \bar{L}}. \quad (57)$$

Therefore, PCDM has an even more dramatic theoretical advantage compared to DQAM in the case when the maximum Lipschitz constant L' is much larger than the average \bar{L} .

6.4 Optimal number of block updates

In this section we propose a simplified model of parallel computing and in it study the performance of a family of parallel coordinate descent methods parameterized by a single parameter: the number of blocks being updated in a single iteration.

In particular, consider the family of PCDMs where S_k is a τ -nice sampling and $\tau \in \{1, 2, \dots, n\}$. Now assume we have $p \in \{1, 2, \dots, n\}$ processors/threads available, each able to compute and apply to the current iterate the update $h^{(i)}(x_k)$ for a single block i , in a unit of time. PCDM, as analyzed, is a synchronous method. That is, a new parallel iteration can only start once the previous one is finished, and hence updating τ blocks will take $\lceil \frac{\tau}{p} \rceil$ amount of time. On the other hand, the iteration complexity of PCDM is better for higher τ . Indeed, by Theorem 6, f satisfies an ESO with respect to \hat{S} with parameters $w = L = (L_1, \dots, L_n)$ and $\beta = \beta(\tau) = 1 + \frac{(\omega-1)(\tau-1)}{n-1}$, where ω is degree of partial separability of f (we assume $n > 1$). If, moreover, $\mu_F(L) = \mu_f(L)$, which is often the case as Ψ is often not strongly convex, then Theorem 10 says that PCDM needs $\frac{n}{\tau}\beta(\tau)c$ iterations, where c is a constant independent of τ , to solve (4) with high probability. Hence, the total amount of time needed for PCDM to solve the problem is equal to

$$T(\tau) = \lceil \frac{\tau}{p} \rceil \frac{n}{\tau} \beta(\tau) c.$$

We can now ask the following natural question: what $\tau \in \{1, 2, \dots, n\}$ minimizes $T(\tau)$? We now show that the answer is $\tau = p$.

Theorem 13. *Assume $f : \mathbf{R}^N \rightarrow \mathbf{R}$ is convex, partially separable of degree ω , and has block Lipschitz gradient with constants $L_1, L_2, \dots, L_n > 0$, where $n > 1$. Further assume $\mu_F(L) = \mu_f(L) > 0$ and consider the family of parallel coordinate descent methods with τ -nice sampling, where $\tau \in \{1, 2, \dots, n\}$, applied to problem (4). Under the parallel computing model with $p \in \{1, 2, \dots, n\}$ processors described above, the method with $\tau = p$ is optimal.*

Proof. We only need to show that $p = \arg \min\{T(\tau) : \tau = 1, 2, \dots, n\}$. It is easy to see that $\frac{n}{\tau}\beta(\tau)$ is decreasing in τ . Since $\lceil \frac{\tau}{p} \rceil$ is constant for $kp + 1 \leq \tau \leq kp$, it suffices to consider $\tau = kp$ for $k = 1, 2, \dots$ only. Finally, $T(kp) = \frac{n}{p}\beta(kp)c$ is increasing in k since $\beta(\cdot)$ is increasing, and we conclude that $k = 1$ and hence $\tau = p$ is optimal. \square

7 Numerical Results

In this section we present two numerical experiments that support the findings of this paper. In both experiments we choose $f(x) = \frac{1}{2}\|b - Ax\|^2$ and $\Psi \equiv 0$.

The first experiment considers the performance of DQAM and the fully parallel variant of PCDM in the above setting where we know that the two methods coincide up to the selection of the stepsize parameters ω and θ (recall Section 5.4). Here we focus on comparing the effects of using the DQAM stepsize $\theta = 1/(2(\omega - 1))$ versus the larger PCDM stepsize $\theta = 1/\omega$.

The second experiment compares DQAM, fully parallel variant of PCDM (i.e., PCDM used with n -nice sampling) and PCDM used with τ -nice sampling, in the situation when the number of available processors is τ , while varying ω (degree of partial separability of f) and τ .

7.1 Impact of the different stepsizes of DQAM and PCDM

Suppose that A has primal block angular structure, $A = [C^T, D^T]^T$, where C is block diagonal ($C \stackrel{\text{def}}{=} \text{diag}(C_1, \dots, C_n)$, for $C_i \in \mathbf{R}^{m_i \times N_i}$), and $D \in \mathbf{R}^{\ell \times N}$ (where $D \stackrel{\text{def}}{=} [D_1, \dots, D_n]$ for $D_i \in \mathbf{R}^{\ell \times N_i}$). Notice that when $D = 0$, the problem is partially separable of degree $\omega = 1$ (i.e., it is fully separable) with respect to the natural block structure (i.e., blocks corresponding to the column submatrices $[C_i; 0; D_i]$). If D is completely dense, the problem is nonseparable ($\omega = n$). In general, ω is equal to the number of matrices D_i that contain at least one nonzero entry.

In this (small scale) experiment we set $n = 100$ and let C_1, \dots, C_{100} be 10% dense matrices of size 150×100 . Subsequently, A is a $15,001 \times 10,000$ sparse matrix. The degree of separability of f varies, and is controlled by setting a subset of the matrices D_1, \dots, D_n to zero.

Twenty five random pairs (A, b) were generated for each $\omega \in \{2, 4, 8, 16, 32, 64, 128\}$, and DQAM and fully parallel PCDM were applied to each problem instance, with a stopping condition of $f(x) - f^* \leq 10^{-4}$. The results of this experiment are presented in Figure 7.1, and all data points are averages over 25 runs.

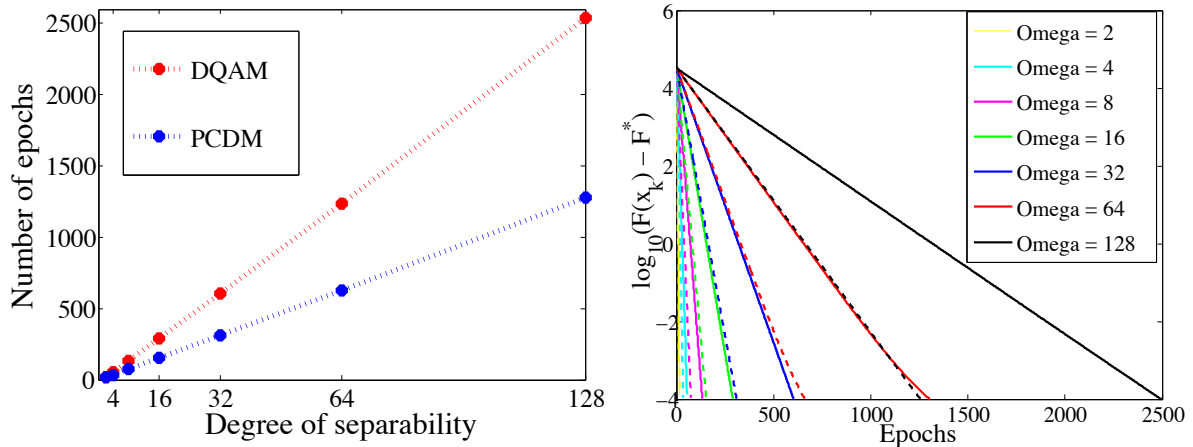


Figure 1: The left plot shows the number of epochs (a full sweep through the data, i.e., all $i = 1, \dots, n$ blocks of x are updated in one epoch) needed to solve the problem as a function of the degree of separability ω . PCDM requires far fewer epochs to reach the stopping condition compared with DQAM. The right plot shows the number of epochs versus the decrease in the function value. The colours correspond to different values of ω , with dotted lines corresponding to PCDM and solid lines corresponding to DQAM. Again, this plot shows the superiority of PCDM and also shows the linear convergence behaviour of the algorithms.

In the left plot in Figure 7.1, we see that when $\omega = 2$, DQAM and PCDM require the same

number of epochs to solve the problem. This is expected because $\theta = 1/(2(\omega - 1)) = 1/2 = 1/\omega$. (Recall Section 5.4, which shows that the two methods coincide in this special case.) Then, as ω grows, PCDM performs far better than DQAM, requiring almost 50% fewer epochs than DQAM. Again, this is predicted by theory because, as $\omega \rightarrow \infty$, $1/(2(\omega - 1)) \rightarrow 1/2\omega$. Furthermore, the right plot in Figure 7.1 shows the linear convergence of the two algorithms.

7.2 Comparison of full vs partial parallelization

Recall that unlike DQAM, PCDM is able to update τ blocks at each iteration, for any τ in the set $\{1, 2, \dots, n\}$, demonstrating useful flexibility of the algorithm. By $\text{PCDM}(\tau)$ we denote the variant of PCDM in which τ blocks are updated at each iteration, using a τ -nice sampling. In this experiment we investigate the performance of DQAM, $\text{PCDM}(n)$ (which in the plots we refer to simply as PCDM) and $\text{PCDM}(\tau)$, for a selection of parameters τ (the number of processors), and ω (the degree of partial separability).

Let us call the time taken for all τ processors to update a single block, one “time unit”. Then, after one time unit of $\text{PCDM}(\tau)$, new gradient information is available to be utilized during the next time unit, which is much earlier than if all n blocks need to be updated in each iteration. On the other hand, for DQAM and PCDM, one iteration corresponds to all n blocks of x being updated. Subsequently, if there are τ processors available, one iteration of DQAM or PCDM (one epoch) corresponds to $\lceil \frac{n}{\tau} \rceil$ time units. However, $\text{PCDM}(\tau)$ will need to perform more iterations than both DQAM and PCDM. When both of these factors are taken into account, we have shown in Theorem 13 that $\text{PCDM}(\tau)$ is optimal in terms of overall complexity if there are τ processors.

The purpose of this experiment is to investigate this phenomenon numerically. Further, let A be a $2 \cdot 10^4 \times 10^4$ sparse matrix, with at most ω nonzero entries per row. Let the stopping condition be $f(x) \leq 10^{-4} b^T b$. The experiment was run for three instances: $\omega = 20, 60, 100$, and for each ω and varying τ , the average number of time units required by DQAM, PCDM and $\text{PCDM}(\tau)$ were recorded. The results are shown in Figure 2.

The colors in Figure 2 correspond to different values of τ . The solid lines correspond to $\text{PCDM}(\tau)$, while the dotted line (respectively dashed line) corresponds to DQAM (respectively PCDM) run with τ processors available. As ω increases, all algorithms require a higher number of time units. Further, as the number of available processors increases, the number of time units decreases. More importantly, for any fixed τ , $\text{PCDM}(\tau)$, requires far fewer time units than PCDM, and both require many fewer time units than DQAM. (Notice the log scale.) This demonstrates the practical advantage of ‘optimizing’ $\text{PCDM}(\tau)$ to the number of available processors, (Section 6.4).

We have also recorded the average cpu time, and the resulting curves are visually indistinguishable from those in Figure 2; only the scale of the vertical axis changes.

7.3 A stochastic optimization example

In this numerical experiment we compare DQAM and PCDM on a stochastic multistage portfolio optimization problem. The problem is, given an initial budget of $S_0 = \$10,000$ to invest in a pool of J assets at time $t = 0$, select assets for inclusion in the portfolio, in such a way that the (expected) profit is maximized at the final timestep $t = T > 0$.

A three-stage portfolio optimization problem ($T = 2$) was generated using real-world data collected from Yahoo! Finance⁴. In particular, monthly (adjusted close) price data was collected

⁴<http://finance.yahoo.com/>

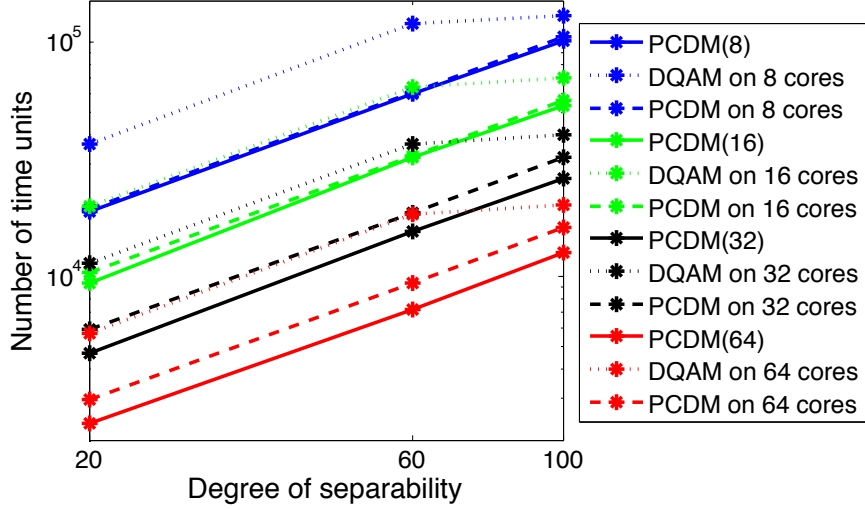


Figure 2: For each fixed $\tau \in \{8, 16, 32, 64\}$, PCDM(τ) (solid line) is better than PCDM (dashed line), and both are far better than DQAM (dotted line).

for 20 assets ($J = 20$) from the FTSE100 for the past two years (September 2011 – August 2013). The price data was used to find the rate-of-return for each asset $j = 1, \dots, J$:

$$r^{(j)}(t_d) = (p^{(j)}(t_d) - p^{(j)}(t_d - 1)) / (p^{(j)}(t_d - 1)), \quad j = 1, \dots, J,$$

where $p^{(j)}(t_d)$ is the price of asset j at the data timestep t_d (i.e., $t_d \in \{\text{September 2011, October 2011, } \dots, \text{August 2013}\}$)⁵. The rate-of-return data for the period September 2011 – August 2012 was used to populate the first time stage of a scenario tree, giving 12 branches from the root node. The remaining data (September 2012 – August 2013) was used to populate the second time stage of the scenario tree, giving 12 branches for *each* of the first stage branches. (A total of $n = 144$ scenarios.) Equal probabilities were assigned to each scenario.

For this portfolio optimization problem, nonanticipativity constraints $b = Ax$ were introduced to link together variables from different scenarios. (See [21] for a thorough discussion of nonanticipativity constraints.) For this application, the degree of separability is $\omega = 2$, and we also set $B_i = I$ and $L_i = 2r$ for $i = 1, \dots, n$. Further, $X_i = \{x^{(i)} \in \mathbf{R}^{160} | Q_i x^{(i)} = q^{(i)}, x^{(i)} \geq 0\}$, where $Q_i x^{(i)} = q^{(i)}$ enforces the inventory and cash-balance equations for the portfolio optimization problem.

The objective function is $g^{(i)}(x^{(i)}) \stackrel{\text{def}}{=} (c^{(i)})^T x^{(i)} = (1 - \gamma) \sum_{i=1}^J \text{prob}^{(i)}(T) \cdot p^{(i)}(0) \cdot x^{(i)}(T)$, where $\gamma = 0.02$ is the transaction cost, $\text{prob}^{(i)}(T)$ is the probability for scenario i at time T , $p^{(i)}(0)$ is the initial price for asset i and $x^{(i)}(T)$ is the amount of asset i held at time T . Therefore, this problem can be written as

$$\max_{x \in X} \sum_{i=1}^n (c^{(i)})^T x^{(i)}, \text{ subject to } Ax = b,$$

⁵The rate-of-return gives the relative increase/decrease in the asset price from one data timestep to the next.

and the augmented Lagrangian is

$$\min_{x \in X} \sum_{i=1}^n -\langle c^{(i)}, x^{(i)} \rangle + \pi^T(b - Ax) + \frac{r}{2} \|b - Ax\|_2^2.$$

In particular, $f = \frac{r}{2} \|b - Ax\|_2^2$ and $\Psi = \sum_{i=1}^n (-\langle c^{(i)}, x^{(i)} \rangle + \langle \pi^T, b - Ax \rangle)$.

The parameters used in the algorithm are, $\theta = 0.5$ (the DQAM stepsize), $\tau = 12$ (the number of processors available), and the convergence stopping test was $\|b - Ax\|_2^2 < 10^{-4}$. The subproblems in these experiments ((11a) and (26a)) were solved using Matlab’s ‘quadprog’ function.

Iterations	DQAM	PCDM-FP	PCDM(τ)
Outer	329	332	292
Total	160,522	158,770	142,710

Table 1: This table reports the number of iterations required by DQAM, fully parallel PCDM, and PCDM optimized to τ processors, for the portfolio optimization problem described in Section 7.3. The iteration counts are broken into, ‘outer’, which gives the number of full Method of Multiplier iterations performed, and ‘total’, which gives the total number of block updates performed throughout the algorithm.

We see that all three algorithms require a similar number of ‘outer’ iterations, with PCDM optimized to $\tau = 12$ processors requiring the least amount of block updates overall, while fully parallel PCDM and DQAM perform similarly.

References

- [1] Arno J. Berger, John M. Mulvey, and Andrzej Ruszczyński. An extension of the DQA algorithm to convex stochastic programs. *SIAM Journal on Optimization*, 4(4):735–753, November 1994.
- [2] D. P. Bertsekas. Extended monotropic programming and duality. *Extended monotropic programming and duality*, 139(2):209–225, 2008.
- [3] Dimitri Bertsekas. *Constrained Optimization and Lagrange Multiplier Methods*. Athena Scientific, 1996.
- [4] E. G. Birgin and J. M. Martínez. Improving ultimate convergence of an augmented Lagrangian method. *Optimization Methods and Software*, 23(2):177–195, 2011.
- [5] Mathieu Blondel, Kazuhiro Seki, and Kuniaki Uehara. Block coordinate descent algorithms for large-scale sparse multiclass classification. *Machine Learning*, 93(1):31–52, October 2013.
- [6] Joseph K. Bradley, Aapo Kyrola, Danny Bickson, and Carlos Guestrin. Parallel coordinate descent for L1-regularized loss minimization. In Lise Getoor and Tobias Scheffer, editors, *28th International Conference on Machine Learning (ICML-11)*, pages 321–328. ACM, June 2011.
- [7] J. Eckstein. Augmented Lagrangian and alternating direction methods for convex optimization: A tutorial and some illustrative computational results. Technical report, December 2012. RUTCOR Research Report RRR 32-2012.

- [8] J. Eckstein and D. P. Bertsekas. On the Douglas-Rachford splitting method and the proximal point algorithm for maximal monotone operators. *Mathematical Programming*, 55:293–318, 1992.
- [9] Olivier Fercoq. Parallel coordinate descent for the Adaboost problem. Technical report, University of Edinburgh, July 2013. arXiv:1310.1840v1 [cs.LG].
- [10] Olivier Fercoq and Peter Richtárik. Smoothed parallel coordinate descent method. Technical report, University of Edinburgh, 2013.
- [11] Adel Hamdi and Andreas Griewank. Properties of an augmented Lagrangian for design optimization. *Optimization Methods and Software*, 25(4):645–664, 2009.
- [12] B. S. He, M. Tao, and X. M. Yuan. Alternating direction method with gaussian back substitution for separable convex programming. *SIAM Journal on Optimization*, 22(2):313–340, 2012.
- [13] Magnus R. Hestenes. Multiplier and gradient methods. *Journal of Optimization Theory and Applications*, 4:303–320, 1969.
- [14] M. Hong and Z. Q. Luo. On the linear convergence of the alternating direction method of multipliers. Technical report, August 2012. arXiv:1208.3922.
- [15] Cho-Jui Hsieh, Kai-Wei Chang, Chih-Jen Lin, S Sathiya Keerthi, and S Sundararajan. A dual coordinate descent method for large-scale linear svm. In *ICML 2008*, pages 408–415, 2008.
- [16] Yin Tat Lee and Aaron Sidford. Efficient accelerated coordinate descent methods and faster algorithms for solving linear systems. Technical report, 2013. arXiv:1305.1922v1.
- [17] Yingying Li and Stanley Osher. Coordinate descent optimization for l_1 minimization with application to compressed sensing; a greedy algorithm. *Inverse Problems and Imaging*, 3:487–503, August 2009.
- [18] Zhaosong Lu and Lin Xiao. On the complexity analysis of randomized block-coordinate descent methods. Technical report, May 2013. arXiv:1305.4723.
- [19] Zhaosong Lu and Lin Xiao. Randomized block coordinate non-monotone gradient method for a class of nonlinear programming. Technical report, June 2013. arXiv:1306.5918.
- [20] John M. Mulvey and Andrzej Ruszczyński. A diagonal quadratic approximation method for large scale linear programs. *Operations Research Letters*, 12:205–215, 1992.
- [21] John M. Mulvey and Andrzej Ruszczyński. A new scenario decomposition method for large scale stochastic optimization. *Operations Research*, 43(3):477–490, 1995.
- [22] Ion Necoara, Yurii Nesterov, and Francois Glineur. Efficiency of randomized coordinate descent methods on optimization problems with linearly coupled constraints. Technical report, June 2012.
- [23] Ion Necoară and A. Patrascu. A random coordinate descent algorithm for optimization problems with composite objective function and linear coupled constraints. *Computational Optimization and Applications*, 57(2):307–337, March 2014.

- [24] Yurii Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM Journal on Optimization*, 22(2):341–362, 2012.
- [25] Andrei Patrascu and Ion Necoara. Efficient random coordinate descent algorithms for large-scale structured nonconvex optimization. Technical report, University Politehnica Bucharest, May 2013.
- [26] Michael J. D. Powell. A method for nonlinear constraints in minimization problems. In Roger Fletcher, editor, *Optimization*, pages 283–298. Academic Press, 1972.
- [27] Zhiwei (Tony) Qin, Katya Scheinberg, and Donald Goldfarb. Efficient block-coordinate descent algorithms for the group lasso. *Mathematical Programming Computation*, 5(2):143–169, June 2013.
- [28] Peter Richtárik and Martin Takáč. Efficient serial and parallel coordinate descent methods for huge-scale truss topology design. In *Operations Research Proceedings 2011*, pages 27–32. Springer, 2012.
- [29] Peter Richtárik and Martin Takáč. Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function. *Mathematical Programming, Ser. A*, 2012.
- [30] Peter Richtárik and Martin Takáč. Parallel coordinate descent methods for big data optimization. Technical report, November 2012. arXiv:1212.0873.
- [31] Peter Richtárik and Martin Takáč. Efficiency of randomized coordinate descent methods on minimization problems with a composite objective function. In *4th Workshop on Signal Processing with Adaptive Sparse Structured Representations*, June 2011.
- [32] R. Tyrell Rockafellar. The multiplier method of Hestenes and Powell applied to convex programming. *Journal of Optimization Theory and Applications*, 12:555–562, 1973.
- [33] R. Tyrell Rockafellar. Augmented Lagrangians and applications of the proximal point algorithm in convex programming. *Mathematics of Operations Research*, 1:97–116, 1976.
- [34] R. Tyrell Rockafellar and Roger J.-B. Wets. Scenarios and policy aggregation in optimization under uncertainty. *Mathematics of Operations Research*, 16:1–23, 1991.
- [35] Andrzej Ruszczyński. An augmented Lagrangian method for block diagonal linear programming problems. *Operations Research Letters*, 8:287–294, 1989.
- [36] Andrzej Ruszczyński. On convergence of an augmented Lagrangian decomposition method for sparse convex optimization. *Mathematics of Operations Research*, 20(3):634–656, 1995.
- [37] Hermann Schwarz. Über einen Grenzübergang durch alternierendes Verfahren. *Vierteljahrsschrift der Naturforschenden Gesellschaft in Zürich*, 15:272–286, 1870.
- [38] Shai Shalev-Shwartz and Ambuj Tewari. Stochastic methods for ℓ_1 regularized loss minimization. In Léon Bottou and Michael Littman, editors, *Proceedings of the 26th International Conference on Machine Learning*, pages 929–936, Montreal, June 2009. Omnipress.

- [39] Shai Shalev-Shwartz and Tong Zhang. Stochastic dual coordinate ascent methods for regularized loss minimization. *Journal of Machine Learning Research*, 14:567–599, 2013.
- [40] George Stephanopoulos and Arthur W. Westerberg. The use of Hestenes’ method of multipliers to resolve dual gaps in engineering system optimization. *Journal of Optimization Theory and Applications*, 15:285–309, 1975.
- [41] Martin Takáč, Avleen Bijral, Peter Richtárik, and Nathan Srebro. Mini-batch primal and dual methods for SVMs. In Sanjoy Dasgupta and David McAllester, editors, *30th International Conference on Machine Learning*, volume 28 of *JMLR: Workshop and Conference Proceedings*, pages 1022–1030, 2013.
- [42] Rachael Tappenden, Peter Richtárik, and Jacek Gondzio. Inexact coordinate descent: complexity and preconditioning. Technical report, University of Edinburgh, April 2013. arXiv:1304.5530.
- [43] Paul Tseng. Convergence of a block coordinate descent method for nondifferentiable minimization. *Journal of Optimization Theory and Applications*, 109:475–494, June 2001.
- [44] X. Wang, M. Hong, S. Ma, and Z. Q. Luo. Solving multiple-block separable convex minimization problems using two-block alternating direction method of multipliers. Technical report, August 2013. arXiv:1308.5294.
- [45] N. Watanabe, Y. Nishimura, and M. Matsubara. Decomposition in large system optimization using the method of multipliers. *Journal of Optimization Theory and Applications*, 25:181–193, 1978.

A Notation Dictionary

For the reader interested in comparing our work with the paper [36] directly, we have included a brief dictionary translating some of the key notation

Ruszczynski [36]	This paper
L	n
N	$\omega - 1$
x_i	$x^{(i)}$
\tilde{x}	x
x	y
$x - \tilde{x}$	$h = y - x$
τ	θ
ρ	r
$\rho\alpha^2$	L'
γ	$\mu_F(e)/2$
$\frac{1}{2}r\ b - \sum_{i=1}^n A_i x_i\ _2^2$	$f(x)$
$f_i(x_i) - \langle A_i^T \pi, x_i \rangle$	$\Psi_i(x^{(i)}) (= g_i(x_i) - \langle A_i^T \pi, x_i \rangle)$
$\Lambda(x)$	$F(x) = f(x) + \Psi(x)$
$\Lambda_i(x_i, \tilde{x})$	$f(x + U_i h^{(i)}) + \Psi_i(y^{(i)})$
$\tilde{\Lambda}(x, \tilde{x})$	$f(x) + \sum_{i=1}^n [f(x + U_i h^{(i)}) - f(x)] + \Psi(x + h)$

B Two Measures of Separability

Here (for the reader's convenience) we confirm the assertion that “ ω_R is the maximum number of blocks linked by any single constraint, decremented by one” made in Section 4.4 and in [36].

We now define the separability measure ω_R introduced for the convex quadratic $f(x) = \frac{r}{2}\|b - Ax\|^2$ by Ruszczyński [36, Section 3] (and called the “number of neighbors” therein).

Let A_{ji} be the j -th row of matrix A_i . Let m_i be the number of nonzero rows in A_i and for each i define an $m \times m_i$ matrix E^i as follows: $E_{jl}^i = 1$ if A_{ji} is the l -th consecutive nonzero row of the matrix A_i , and 0 otherwise. Note that E^i is a matrix containing zeros and m_i ones, one in each column. Let E_u^i be the u -th column of matrix E^i . Then, for any $i \in \{1, 2, \dots, n\}$ and $u \in \{1, 2, \dots, m_i\}$ define

$$V(i, u) \stackrel{\text{def}}{=} \{(i', u') : i' \in \{1, 2, \dots, n\}, u' \in \{1, 2, \dots, m_{i'}\}, i' \neq i, \langle E_u^i, E_{u'}^{i'} \rangle \neq 0\}. \quad (58)$$

Definition 14. *The Ruszczyński degree of separability of the function f defined in (5) is*

$$\omega_R = \max\{|V(i, u)| : i = 1, 2, \dots, n, u = 1, 2, \dots, m_i\}. \quad (59)$$

Using the following arguments, we show that the Ruszczyński degree of separability ω_R is indeed the maximum number of blocks linked by a single constraint decremented by one. Let us fix $i \in \{1, 2, \dots, n\}$, $u \in \{1, 2, \dots, m_i\}$ and let $j = j(i, u)$ be a row index for which $E_{ju}^i = 1$. Since E^i is a 0-1 matrix with exactly one entry of each column equal to 1, we have $E_{j'u}^i = 0$ for all $j' \neq j$. This means that for any $i' \in \{1, 2, \dots, n\}$ and $u' \in \{1, 2, \dots, m_{i'}\}$,

$$\langle E_u^i, E_{u'}^{i'} \rangle \neq 0 \Leftrightarrow E_{ju'}^{i'} = 1. \quad (60)$$

Likewise, $E^{i'}$ has at most one entry equal to one in each row. Moreover, the j -th row of $E^{i'}$ contains 1 precisely when $A_{ji'} \neq 0$. This means that $|\{u' : E_{ju'}^{i'} = 1\}|$ is 1 if $A_{ji'} \neq 0$ and 0 otherwise. So

$$\begin{aligned} |V(i, u)| &\stackrel{(58)}{=} |\{(i', u') : i' \neq i, \langle E_u^i, E_{u'}^{i'} \rangle \neq 0\}| \\ &\stackrel{(60)}{=} |\{(i', u') : i' \neq i, E_{ju'}^{i'} = 1\}| \\ &= \sum_{i' \neq i} |\{u' : E_{ju'}^{i'} = 1\}| \stackrel{(31)}{=} \omega_j - 1. \end{aligned} \quad (61)$$

Building on this result we can now write

$$\begin{aligned} \omega_R &\stackrel{(59)}{=} \max\{|V(i, u)| : i \in \{1, 2, \dots, n\}, u \in \{1, 2, \dots, m_i\}\} \\ &\stackrel{(61)}{=} \max\{\omega_{j(i, u)} - 1 : i \in \{1, 2, \dots, n\}, u \in \{1, 2, \dots, m_i\}\} \\ &= \max_{j \in \{1, 2, \dots, m\}} \omega_j - 1 \stackrel{(32)}{=} \omega - 1. \end{aligned}$$

In the third identity above we used the simple observation that every row $j \in \{1, 2, \dots, m\}$ for which $\omega_j \neq 0$ can be written as $j = j(i, u)$ for any i for which $A_{ji} \neq 0$, and some u (which depends on i).

Indeed, ω_R is the maximum number of blocks linked by any single constraint, minus one.