

Towards Totally Asynchronous Primal-Dual Convex Optimization in Blocks

Katherine R. Hendrickson and Matthew T. Hale*

Abstract

We present a parallelized primal-dual algorithm for solving constrained convex optimization problems. The algorithm is “block-based,” in that vectors of primal and dual variables are partitioned into blocks, each of which is updated only by a single processor. We consider four possible forms of asynchrony: in updates to primal variables, updates to dual variables, communications of primal variables, and communications of dual variables. We explicitly construct a family of counterexamples to rule out permitting asynchronous communication of dual variables, though the other forms of asynchrony are permitted, all without requiring bounds on delays. A first-order update law is developed and shown to be robust to asynchrony. We then derive convergence rates to a Lagrangian saddle point in terms of the operations agents execute, without specifying any timing or pattern with which they must be executed. These convergence rates contain a synchronous algorithm as a special case and are used to quantify an “asynchrony penalty.” Numerical results illustrate these developments.

I. INTRODUCTION

A wide variety of machine learning problems can be formalized as convex programs [3], [6], [20], [21]. Large-scale machine learning then requires solutions to large-scale convex programs, which can be accelerated through parallelized solvers running on networks of processors. In large networks, it is difficult (or outright impossible) to synchronize their behaviors. The behaviors of interest are computations, which generate new information, and communications, which share this new information with other processors. Accordingly, we are interested in asynchrony-tolerant large-scale optimization.

The challenge of asynchrony is that it causes disagreements among processors that result from receiving different information at different times. One way to reduce disagreements is through repeated averaging of processors’ iterates. This approach dates back several decades [24], and approaches of this class include primal [16]–[18], dual [7], [22], [23], and primal-dual [14], [26] algorithms. However, these averaging-based methods require bounded delays in some form, often through requiring connectedness of agents’ communication graphs over intervals of a prescribed length [4, Chapter 7]. In some applications, delays are outside agents’ control, e.g., in a contested environment where communications are jammed, and thus delay bounds cannot be reliably enforced. Moreover, graph connectivity cannot be easily checked locally by individual agents, meaning even satisfaction or violation of connectivity bounds is not readily ascertained. In addition, these methods require multiple processors to update each decision variable,

This work was supported by a task order contract from the Air Force Research Laboratory through Eglin AFB.

*The authors are with the Department of Mechanical and Aerospace Engineering, Herbert Wertheim College of Engineering, University of Florida. Emails: {kat.hendrickson, matthewhale}@ufl.edu.

which duplicates computations and increases processors' workloads. This can be prohibitive in large problems, such as learning problems with billions of data points.

Therefore, in this paper we develop a parallelized primal-dual method for solving large constrained convex optimization problems. Here, by "parallelized," we mean that each decision variable is updated only by a single processor. As problems grow, this has the advantage of keeping each processor's computational burden approximately constant. The decision variables assigned to each processor are referred to as a "block," and block-based algorithms date back several decades as well [2], [24]. Those early works solve unconstrained or set-constrained problems, in addition to select problems with functional constraints [4]. To bring parallelization to arbitrary constrained problems, we develop a primal-dual approach that does not require constraints to have a specific form.

Block-based methods have previously been shown to tolerate arbitrarily long delays in both communications and computations in unconstrained problems [2], [13], [25], eliminating the need to enforce and verify delay boundedness assumptions. For constrained problems of a general form, block-based methods have been paired with primal-dual algorithms with centralized dual updates [10], [12] and/or synchronous primal-updates [15]. To the best of our knowledge, arbitrarily asynchronous block-based updates have not been developed for convex programs of a general form. A counterexample in [12] suggested that arbitrarily asynchronous communications of dual variables can preclude convergence, though that example leaves open the extent to which dual asynchrony is compatible with convergence.

In this paper, we present a primal-dual optimization algorithm that permits arbitrary asynchrony in primal variables, while accommodating dual asynchrony to the extent possible. Four types of asynchrony are possible: (i) asynchrony in primal computations, (ii) asynchrony in communicating primal variables, (iii) asynchrony in dual computations, (iv) asynchrony in communicating dual variables. The first contribution of this paper is to show that item (iv) is fundamentally problematic using an explicit family of counterexamples that we construct. This family shows, in a precise way, that even small disagreements among dual variables can cause primal computations to diverge. For this reason, we rule out asynchrony in communicating dual variables. However, we permit all other forms of asynchrony, and, relative to existing work, this is the first algorithm to permit arbitrarily asynchronous computations of dual variables in blocks.

The second contribution of this paper is to establish convergence rates. These rates are shown to depend upon problem parameters, which lets us calibrate their values to improve convergence. Moreover, we show that convergence can be inexact due to dual asynchrony, and thus the scalability of parallelization comes at the expense of a potentially inexact solution. We term this inexactness the "asynchrony penalty," and we give an explicit bound on it. Simulation results show convergence of this algorithm in practice, and illustrate concretely that the asynchrony penalty is slight.

The rest of the paper is organized as follows. Section II provides the necessary background on convex optimization and formally gives the asynchronous primal-dual problem statement. Then Section III discusses four possible asynchronous behaviors, provides a counterexample to complete asynchrony, and presents our asynchronous algorithm. Primal and dual convergence rates are developed in Section IV. Section V presents a numerical example with

implications for relationships among parameters. Finally, we present our conclusions in Section VI.

II. BACKGROUND AND PROBLEM STATEMENT

We study the following form of optimization problem.

Problem 1: Given $h : \mathbb{R}^n \rightarrow \mathbb{R}$, $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$, and $X \subset \mathbb{R}^n$, asynchronously solve

$$\begin{aligned} & \text{minimize } h(x) \\ & \text{subject to } g(x) \leq 0 \\ & \qquad \qquad x \in X. \end{aligned} \quad \diamond$$

We assume the following about the objective function h .

Assumption 1: h is twice continuously differentiable and convex. △

We make a similar assumption about the constraints g .

Assumption 2: g satisfies Slater's condition, i.e., there exists $\bar{x} \in X$ such that $g(\bar{x}) < 0$. For all $j \in \{1, \dots, m\}$, the function g_j is twice continuously differentiable and convex. △

Assumptions 1 and 2 permit a wide range of functions to be used, such as all convex polynomials of all orders.

We impose the following assumption on the constraint set.

Assumption 3: X is non-empty, compact, and convex. It can be decomposed into $X = X_1 \times \dots \times X_N$. △

Assumption 3 permits many sets to be used, such as box constraints, which often arise multi-agent optimization [19].

We will solve Problem 1 using a primal-dual approach. This allows the problem to be parallelized across many processors by re-encoding constraints through Karush-Kuhn-Tucker (KKT) multipliers. In particular, because the constraints g couple the processors' computations, they can be difficult to enforce in a distributed way. However, by introducing KKT multipliers to re-encode constraints, we can solve an equivalent, higher-dimensional unconstrained problem.

An ordinary primal-dual approach would find a saddle point of the Lagrangian associated with Problem 1, defined as $L(x, \mu) = h(x) + \mu^T g(x)$. That is, one would solve $\min_x \max_\mu L(x, \mu)$. However, L is affine in μ , which implies that $L(x, \cdot)$ is concave but not strongly concave. Strong convexity has been shown to provide robustness to asynchrony in minimization problems [4], and thus we wish to endow the maximization over μ with strong concavity. We use a Tikhonov regularization [8, Chapter 12] in μ to form

$$L_\delta(x, \mu) = h(x) + \mu^T g(x) - \frac{\delta}{2} \|\mu\|^2, \quad (1)$$

where $\delta > 0$.

Instead of regularizing with respect to the primal variable x , we impose the following assumption in terms of the Hessian $H(x, \mu) = \nabla_x^2 L_\delta(x, \mu)$. When convenient, we suppress the arguments x and μ .

Assumption 4 (Diagonal Dominance): The $n \times n$ Hessian matrix $H = \nabla_x^2 L_\delta(x, \mu)$ is β -diagonally dominant. That is, for all i from $1, \dots, n$,

$$|H_{ii}| - \beta \geq \sum_{\substack{j=1 \\ j \neq i}}^n |H_{ij}|. \quad \triangle$$

If this assumption does not hold, the Lagrangian can be regularized with respect to the primal variable as well, leading to H 's diagonal dominance. Some problems satisfy this assumption without regularizing [9], and, for such problems, we proceed without regularizing to avoid regularization error.

Under Assumptions 1-4, Problem 1 is equivalent to the following saddle point problem.

Problem 2: Let Assumptions 1-4 hold and fix $\delta > 0$. For L_δ defined in Equation (1), asynchronously compute

$$(\hat{x}_\delta, \hat{\mu}_\delta) := \arg \min_{x \in X} \arg \max_{\mu \in \mathbb{R}_+^m} L_\delta(x, \mu). \quad \diamond$$

It is in this form that we will solve Problem 1, and we present our algorithm for doing so in the next section.

III. ASYNCHRONOUS PRIMAL-DUAL ALGORITHM

One challenge of Problem 2 is that $\hat{\mu}_\delta$ is maximized over the unbounded domain \mathbb{R}_+^m , which is the non-negative orthant of \mathbb{R}^m . Because this domain is unbounded, gradients and other terms are unbounded, which makes convergence analysis challenging as dual iterates may not be within a finite distance of the optimum. To remedy this problem, we next compute a non-empty, compact, and convex set M that contains $\hat{\mu}_\delta$.

Lemma 1: Let Assumptions 1-4 hold, let \bar{x} be a Slater point of g , and set $h^* := \min_{x \in X} h(x)$. Then

$$\hat{\mu}_\delta \in M := \left\{ \mu \in \mathbb{R}_+^m : \|\mu\|_1 \leq \frac{h(\bar{x}) - h^*}{\min_{1 \leq j \leq m} -g_j(\bar{x})} \right\}.$$

Proof: Follows Section II-C in [11]. ■

Here, h^* denotes the optimal unconstrained objective function value, though any lower-bound for this value will also provide a valid M . In particular, h is often non-negative and one can substitute 0 in place of h^* in such cases.

Solving Problem 2 asynchronously requires choosing an update law that we expect to be robust to asynchrony and simple to implement in a distributed fashion. In this context, first-order gradient-based methods offer both some degree of inherent robustness, as well as computations that are simpler than other available methods, such as Newton-type methods. We apply a projected gradient method to both the primal and dual variables, which is shown in Algorithm 1, and is based on the seminal Uzawa iteration [1].

Algorithm 1

Let $x(0)$ and $\mu(0)$ be given. For values $k = 0, 1, \dots$, execute

$$\begin{aligned} x(k+1) &= \Pi_X[x(k) - \gamma(\nabla_x L_\delta(x(k), \mu(k)))] \\ \mu(k+1) &= \Pi_M[\mu(k) + \gamma(\nabla_\mu L_\delta(x(k), \mu(k)))] \end{aligned} \quad (2)$$

where γ is a step-size, Π_X is the Euclidean projection onto X , and Π_M is the Euclidean projection onto M .

A. Overview of Approach

We are interested in distributing Algorithm 1 among a number of processors while allowing agents to generate and share information as asynchronously as possible. We consider N agents indexed over $a \in [N] := \{1, \dots, N\}$. We partition the set $[N]$ into $[N_p]$ and $[N_d]$ (where $N_p + N_d = N$). The set $[N_p]$ contains indices of agents that update primal variables (contained in x), while $[N_d]$ contains indices of agents that update dual variables (contained in μ). Using a primal-dual approach, there are four behaviors that could be asynchronous: (i) computations of updates to primal variables, (ii) communications to share updated values of primal variables, (iii) computations of updates to dual variables, and (iv) communications to share updated values of dual variables.

(i) *Computations of Updates to Primal Variables:* When parallelizing Equation (2) across the N_p primal agents, we index all primal agents' computations using the same iteration counter, $k \in \mathbb{N}$. However, they may perform updates at different times. The subset of times at which primal agent $i \in [N_p]$ computes an update is denoted by $K^i \subset \mathbb{N}$. For distinct $i, j \in [N_p]$, K^i and K^j need not have any relationship.

(ii) *Communications of Updated Primal Variables:* Primal variable communications are also totally asynchronous. A primal variable's current value may be sent to other primal and dual agents that need it at each time k . We use the notation P_j^i to denote the set of times¹ at which primal agent i sends values of its primal variables to agent j . Similarly, we use the notation D_c^i to denote the set of times at which primal agent i sends updated values to dual agent $c \in [N_d]$.

(iii) *Computations of Updates to Dual Variables:* Dual agents wait for each primal agent's updated state before computing an update. Dual agents may perform updates at different times because they may receive primal updates at different times. In some cases, a dual agent may receive multiple updates from a subset of primal agents prior to receiving all required primal updates. In this case, only the most recently received update from a primal agent will be used in the dual agent's computation. For all $c \in [N_d]$, dual agent c keeps an iteration count t_c to track the number of updates it has completed.

(iv) *Communications of Updated Dual Variables:* Previous work [12] has shown that allowing primal agents to disagree arbitrarily about dual variables can outright preclude convergence. This is explained by the following: fix $\mu^1, \mu^2 \in M$. Then an agent with μ^1 onboard is minimizing $L(\cdot, \mu^1)$, while an agent with μ^2 onboard is minimizing $L(\cdot, \mu^2)$. If μ^1 and μ^2 are arbitrarily far apart, then it is not surprising that the minima of $L(\cdot, \mu^1)$ and $L(\cdot, \mu^2)$ are as well. However, one may conjecture that small disagreements in dual variables lead to small distances between these minima. Below, we show that this conjecture is false and that even small disagreements in dual variables can lead to arbitrarily large distances between the minima they induce. Even limited asynchrony can lead to small disagreements in dual variables, and, in light of the above discussion, this can cause primal agents' computations to reach points that are arbitrarily far apart.

¹We assume that there is no delay between sending and receiving messages. There is no loss of generality in our results because we can make P_j^i and D_c^i the times at which messages are received. However, assuming zero delays simplifies the forthcoming discussion and analysis and is done for the remainder of the paper.

Therefore, we will develop an algorithm that proceeds with all agents agreeing on the value of μ while still allowing dual computations to be divided among dual agents. This is accomplished by allowing primal agents to work completely asynchronously (updates are computed and sent at different times to different agents) but requiring that dual updates are sent to all primal agents at the same time². After a dual agent computes an update, it sends its updated dual variable to all primal agents. Again, for simplicity, it is assumed that the update is received at the same time it is sent. When dual agent c sends the updated dual variable μ_c^c to all primal agents, it also sends its iteration count t_c . This allows primal agents to annotate which version of μ is used in their updates. Primal agents disregard any received primal updates that use an outdated version of μ_c (as indicated by t_c). This ensures that primal updates are not mixed if they rely on different dual values.

B. Counterexample to the Asynchronous Dual Case

Below we show that behavior (iv) above, communications to share updated values of dual variables, cannot be asynchronous in general. The intuition here is as follows. In a primal-dual setup, one can regard each fixed choice of dual vector as specifying a problem to solve in a parameterized family of minimization problems. Formally, with μ fixed, agents solve

$$\text{minimize}_{x \in X} L_\delta(x, \mu) := h(x) + \mu^T g(x) - \frac{\delta}{2} \|\mu\|^2.$$

For two primal agents with different values of μ , denoted μ^1 and μ^2 , they solve two different problems: agent 1 minimizes $L_\delta(\cdot, \mu^1)$ while agent 2 minimizes $L_\delta(\cdot, \mu^2)$. With a gradient-based method to minimize over x , gradients depend linearly upon μ . This may lead one to believe that for

$$\hat{x}_1 := \arg \min_{x \in X} L_\delta(x, \mu^1) \quad \text{and} \quad \hat{x}_2 := \arg \min_{x \in X} L_\delta(x, \mu^2),$$

having $\|\mu^1 - \mu^2\|$ small implies that $\|\hat{x}_1 - \hat{x}_2\|$ is also small. However, we show in the following theorem that this is false.

Theorem 1: Fix any $\epsilon > 0$ and $L > \epsilon$. Then, under Assumptions 1-4, there always exists a problem such that $\|\mu^1 - \mu^2\| < \epsilon$ and $\|\hat{x}_1 - \hat{x}_2\| > L$.

Proof: See the appendix. ■

C. Glossary of Notation

Every agents stores a local copy of x and μ for use in local computations. The following notation is used in our formal algorithm statement below.

D_c^i The times at which messages are sent by primal agent i and received by dual agent c .

$g_c(x)$ The c^{th} entry of the constraint function, g , evaluated at x .

k The iteration count used by all primal agents.

²Even if they are not sent and/or received at the same time, we can apply any procedure to synchronize these values and the algorithm will remain the same. We assume synchrony in sending and receiving these values merely to simplify the forthcoming discussion.

- K^i The set of times at which primal agent i performs updates.
- \mathcal{N}_i Essential neighborhood of agent i . Agent j is an essential neighbor of agent i if $\nabla_{x_j} L_\delta$ depends upon x_i .
Then agent i communicates with agent j to ensure it has the information necessary to compute gradients.
- $[N_d]$ Set containing the indices of all dual agents.
- $[N_p]$ Set containing the indices of all primal agents.
- P_j^i The times at which messages are sent by primal agent i and received by primal agent j .
- $\tau_j^i(k)$ Time at which primal agent j computed the update that it sent agent i at time k (i can be primal or dual).
Note that $\tau_i^i(k) = k$ for all $i \in [N_p]$.
- t The vector of dual agent iteration counts. The c^{th} entry, t_c , is the iteration count for dual agent c 's updates.
- t_c The iteration count for dual agent c 's updates. This is sent along with μ_c^c to all agents.
- v_i^c The iteration count used by dual agent c for updates received from primal agent i between dual updates.
- x_j^i Agent i 's value for the primal variable j , which is updated/sent by primal agent j . If agent i is primal, it is indexed by both k and t ; if agent i is dual it is indexed by t .
- $x_i^i(k; t)$ Agent i 's value for its primal variable i at primal time k , calculated with dual update t .
- $x^*(t)$ The fixed point of $f = \Pi_X [x - \gamma \nabla_x L_\delta(x, \mu(t))]$ with respect to a fixed $\mu(t)$.
- x_t^c Abbreviation for $x^c(t)$, which is dual agent c 's copy of the primal vector at time t .
- \hat{x}_δ The primal component of the saddle point of L_δ . Part of the optimal solution pair $(\hat{x}_\delta, \hat{\mu}_\delta)$.
- \hat{x}^t Given $\mu(t)$, $\hat{x}^t = \arg \min_{x \in X} L_\delta(x, \mu(t))$.
- μ_d^c Agent c 's copy of dual variable d , which is updated/sent by dual agent d . Agent c may be primal or dual.
- $\hat{\mu}_\delta$ The dual component of the saddle point of L_δ , $(\hat{x}_\delta, \hat{\mu}_\delta)$.
- $\hat{\mu}_{c,\delta}$ The c^{th} entry of $\hat{\mu}_\delta$.
- M_c The set $\{\nu \in \mathbb{R}_+ : \nu \leq \frac{h(\bar{x}) - h^*}{\min_j -g_j(\bar{x})}\}$. This uses the upper bound in Lemma 1 to project individual components of μ .

D. Statement of Algorithm

Having defined our notation, we impose the following assumption on agents' communications and computations.

Assumption 5: For all $i \in [N_p]$, the set K^i is infinite. If $\{k_n\}_{n \in \mathbb{N}}$ is an increasing set of times in K^i , then $\lim_{n \rightarrow \infty} \tau_i^j(k_n) = \infty$ for all $j \in [N_p]$ and $\lim_{n \rightarrow \infty} \tau_i^c(k_n) = \infty$ for all $c \in [N_d]$. \triangle

This simply ensures that no agent ever stop computing or communicating, though delays can be arbitrarily large.

We now define the asynchronous primal-dual algorithm.

Algorithm 2

Step 0: Initialize all primal and dual agents with $x(0) \in X$ and $\mu(0) \in M$. Set $t = 0$ and $k = 0$.

Step 1: For all $i \in [N_p]$ and all $j \in \mathcal{N}_i$, if $k \in P_j^i$, then agent j sends $x_j^i(k; t)$ to agent i .

Step 2: For all $i \in [N_p]$ and all $c \in [N_d]$, if agent i receives a dual variable update from agent c , it uses the accompanying t_c to update the vector t and performs the update

$$\mu_c^i(t) = \mu_c^c(t_c).$$

Step 3: For all $i \in [N_p]$ and all $j \in \mathcal{N}_i$, execute

$$x_i^i(k+1; t) = \begin{cases} \Pi_{X_i} [x_i^i(k; t) - \gamma \nabla_{x_i} L_\delta(x_i^i(k; t), \mu(t))] & k \in K^i \\ x_i^i(k; t) & k \notin K^i \end{cases}$$

$$x_j^i(k+1; t) = \begin{cases} x_j^j(\tau_j^i(k+1); t) & i \text{ receives } j\text{'s state at } k+1 \\ x_j^i(k; t) & \text{otherwise} \end{cases}.$$

Step 4: If $k+1 \in D_c^i$, agent i sends $x_i^i(k+1; t)$ to dual agent c . Set $k := k+1$.

Step 5: For $c \in [N_d]$ and $i \in [N_p]$, if dual agent c receives an update from primal agent i computed with dual update t , it sets

$$x_i^c(t_c) = \begin{cases} x_i^i(\tau_i^c(k), t) & x_i^i \text{ received prior to update } t_c + 1 \\ x_i^c(t_c - 1) & \text{otherwise} \end{cases}.$$

Step 6: For $c \in [N_d]$, if agent c has received an update from every primal agent for the latest dual iteration t , it executes

$$\mu_c^c(t_c + 1) = \Pi_{M_c} [\mu_c^c(t_c) + \rho \frac{\partial L}{\partial \mu_c}(x^c(t_c), \mu^c(t_c))].$$

Step 7: If dual agent c updated in Step 6, then it sends $\mu_c^c(t_c + 1)$ to all primal agents. Set $t_c := t_c + 1$.

Step 8: Return to Step 1.

IV. CONVERGENCE

To define an overall convergence rate to the optimal solution $(\hat{x}_\delta, \hat{\mu}_\delta)$, we first fix the dual variable μ and find the primal convergence rate. We then find the overall dual convergence rate to $\hat{\mu}_\delta$ by showing that dual variables converge to $\hat{\mu}_\delta$ over time, which lets us show that primal variables converge to \hat{x}_δ .

A. Primal Convergence with Fixed Dual Variable

Given a fixed $\mu(t)$, projected gradient descent for minimizing $L_\delta(\cdot, \mu(t))$ may be written as

$$f(x) = \Pi_X [x - \gamma \nabla_x L_\delta(x, \mu(t))], \quad (3)$$

where $\gamma > 0$. Leveraging some existing theoretical tools in the study of optimization algorithms [2], [5], we can study f in a way that elucidates its behavior under asynchrony.

According to [5], the assumption of diagonal dominance guarantees that f has the contraction property

$$\|f(x) - x^*(t)\| \leq \alpha \|x - x^*(t)\| \text{ for all } x \in X,$$

where $\alpha \in [0, 1)$ and $x^*(t)$ is a fixed point of f , which depends on the choice of fixed $\mu(t)$. However, the value of α is not specified in [5], and it is precisely that value that governs the rate of convergence to a solution. We therefore compute α explicitly below. First, we bound the step-size γ .

Definition 1: Define the primal step-size $\gamma > 0$ to satisfy

$$\gamma < \frac{1}{\max_i \max_{x \in X} \max_{\mu \in M} \sum_{j=1}^n |H_{ij}(x, \mu)|}.$$

Because x and μ both take values in compact sets, each entry H_{ij} is bounded above and below, and thus the upper bound on γ is positive.

Following the method in [2], two $n \times n$ matrices G and F must also be defined.

Definition 2: Define the $n \times n$ matrices G and F as

$$G = \begin{bmatrix} |H_{11}| & -|H_{12}| & \dots & -|H_{1n}| \\ \vdots & \vdots & \ddots & \vdots \\ -|H_{n1}| & -|H_{n2}| & \dots & |H_{nn}| \end{bmatrix} \text{ and } F = I - \gamma G,$$

where I is the $n \times n$ identity matrix.

We now state the following lemma that relies on meeting the conditions listed in [2].

Lemma 2: Let f , G , and F be as above and let Assumptions 1-4 hold. Then, $|f(x) - f(y)| \leq F|x - y|$, for all $x, y \in \mathbb{R}^n$, where $|v|$ denotes the element-wise absolute value of the vector $v \in \mathbb{R}^n$ and the inequality holds component-wise.

Proof: Three conditions must be satisfied in [2]: (i) γ is sufficiently small, (ii) G is positive definite, and (iii) F is positive definite.

(i) *γ is sufficiently small:* Results in [5] require $\gamma \sum_{j=1}^n |H_{ij}| < 1$ for all $i \in \{1, \dots, n\}$, which here follows immediately from Definition 1.

(ii) *G is positive definite:* By definition, G has only positive diagonal entries. By H 's diagonal dominance we have the following inequality for all $i \in \{1, \dots, n\}$:

$$|G_{ii}| = |H_{ii}| \geq \sum_{\substack{j=1 \\ j \neq i}}^n |H_{ij}| + \beta > \sum_{\substack{j=1 \\ j \neq i}}^n |H_{ij}| = \sum_{\substack{j=1 \\ j \neq i}}^n |G_{ij}|.$$

Because G has positive diagonal entries, is symmetric, and is strictly diagonally dominant, G is positive definite by Gershgorin's Circle Theorem.

(iii) *F is positive definite:* Definition 1 ensures the diagonal entries of F are always positive. And F is diagonally dominant if, for all $i \in \{1, \dots, n\}$,

$$|F_{ii}| = 1 - \gamma |H_{ii}| > \gamma \sum_{\substack{j=1 \\ j \neq i}}^n |H_{ij}| = \sum_{\substack{j=1 \\ j \neq i}}^n |F_{ij}|.$$

This can be rewritten as $\gamma \sum_{j=1}^n |H_{ij}| < 1$, which was satisfied under Condition (i). Because F has positive diagonal entries, is symmetric, and is strictly diagonally dominant, F is positive definite by Gershgorin's Circle Theorem. ■

To establish convergence in x , we will show that $\|f(x) - f(x^*(t))\| \leq \alpha \|x - x^*(t)\|$ for all $x \in X$, where $\alpha \in [0, 1)$. Toward doing so, we define the following.

Definition 3: Let $v \in \mathbb{R}^n$. Then $\|v\|_{max} = \max_i |v_i|$.

In this work, we consider a scalar maximum norm because we consider agents that update scalar blocks, though updating non-scalar blocks is readily accommodated by considering a block-maximum norm [4].

We next show that the gradient update law f in Equation (3) converges with asynchronous, distributed computations. Furthermore, we quantify convergence in terms of γ and β .

Lemma 3: Let f, H, G, F, γ , and $\|v\|_{max}$ be as defined above. Let Assumptions 1-4 hold and fix $\mu(t) \in M$. Then for a fixed point $x^*(t)$ of f and for all $x \in \mathbb{R}^n$,

$$\|f(x) - f(x^*(t))\|_{max} \leq q_p \|x - x^*(t)\|_{max},$$

where $q_p = (1 - \gamma\beta) \in [0, 1)$.

Proof: For notational simplicity we write $x^*(t)$ simply as x^* . Assumption 4 and the definition of F give

$$\sum_{j=1}^n F_{ij} = 1 - \gamma \left(|H_{ii}| - \sum_{\substack{j=1 \\ j \neq i}}^n |H_{ij}| \right) \leq 1 - \gamma\beta.$$

This result, Definition 3, and Lemma 2 give

$$\begin{aligned} \|f(x) - f(x^*)\|_{max} &= \max_i |f_i(x) - f_i(x^*)| \\ &\leq \max_i \sum_{j=1}^n F_{ij} |x_j - x_j^*| \leq \max_l |x_l - x_l^*| \max_i \sum_{j=1}^n F_{ij} \\ &\leq \max_l |x_l - x_l^*| (1 - \gamma\beta) = (1 - \gamma\beta) \|x - x^*\|_{max}. \end{aligned}$$

All that remains is to show $(1 - \gamma\beta) \in [0, 1)$. By Definition 1 and the inequality $|H_{ii}| \geq \beta$, for all $x \in X$ and $\mu(t) \in M$,

$$\gamma\beta < \frac{\beta}{\max_i \sum_{j=1}^n |H_{ij}|} \leq \frac{\beta}{\max_i |H_{ii}|} < \frac{\beta}{\beta} = 1. \quad \blacksquare$$

The primal-only convergence rate can be computed by leveraging results in [12] in terms of the number of operations the primal agents have completed (counted in the appropriate sequence). Namely, we count operations as follows. For a given dual variable with iteration vector t , we set $\text{ops}(k, t) = 0$. Then, after all primal agents have computed an update to their decision variable and sent it to and had it received by all other primal agents that need it, say by time k' , we increment ops to $\text{ops}(k', t) = 1$. After $\text{ops}(k', t) = 1$, we then wait until all primal agents have subsequently computed a new update (still using the same dual variable indexed with t) and it has been received by all other primal agents that need it. If this occurs at time k'' , then we set $\text{ops}(k'', t) = 2$, and then this process continues. If at some time k''' , primal agents receive an updated μ (whether just a single dual agent sent

an update or multiple agents send updates) with an iteration vector of t' , then the cycle count would begin again with $\text{ops}(k''', t') = 0$.

Theorem 2: Let Assumptions 1-4 hold. For $\mu(t)$ fixed and the agents asynchronously executing the gradient update law f ,

$$\|x^i(k) - x^*(t)\|_{max} \leq q_p^{\text{ops}(k,t)} \max_j \|x^j(k) - x^*(t)\|_{max},$$

where $x^*(t)$ is the fixed point of f with $\mu(t)$ fixed.

Proof: We write x^* in place of $x^*(t)$ for simplicity. From Lemma 3 we see that f is a q_p -contraction mapping with respect to the norm $\|\cdot\|_{max}$. From Section 6.3 in [5], this contraction property implies that there exist sets of the form

$$X(k) = \{x \in \mathbb{R}^n \mid \|x - x^*\|_{max} \leq q_p^k \|x(0) - x^*\|_{max}\}$$

that satisfy the following criteria from [12]:

- i. $\dots \subset X(k+1) \subset X(k) \subset \dots \subset X$
- ii. $\lim_{k \rightarrow \infty} X(k) = \{x^*\}$
- iii. For all i , there are sets $X_i(k) \subset X_i$ satisfying

$$X(k) = X_1(k) \times \dots \times X_N(k)$$

- iv. For all $y \in X(k)$ and all $i \in [N_p]$, $f_i(y) \in X_i(k+1)$, where $f_i(y) = \Pi_{X_i} [y_i - \gamma \nabla_{x_i} L_\delta(y)]$.

We will use these properties to compute the desired convergence rate. Suppose all agents have a fixed $\mu(t)$ onboard. Upon receipt of this $\mu(t)$, agent i has $x^i(k; t) \in X(0)$ by definition. Suppose at time ℓ_i that agent i computes a state update. Then $x_i^i(\ell_i + 1; t) \in X_i(1)$. For $m = \max_{i \in [N_p]} \ell_i + 1$, we find that $x_i^i(m; t) \in X_i(1)$ for all i . Next, suppose that, after all updates have been computed, these updated values are sent to and received by all agents that need them, say at time m' . Then, for any $i \in [N_p]$, agent i has $x_j^i(m'; t) \in X_j(1)$ for all $j \in [N_p]$. In particular, $x^i(m'; t) \in X(1)$, and this is satisfied precisely when a single cycle has occurred. Iterating this argument for subsequent cycles completes the proof. ■

B. Dual convergence

We next derive a componentwise convergence rate for the dual variable.

Theorem 3: Let Assumptions 1-4 hold. Fix $\delta > 0$ and let $\rho \in \left(\frac{3-\sqrt{3}}{3\delta}, \frac{3+\sqrt{3}}{3\delta}\right)$. Then

$$|\mu_c^c(t_c + 1) - \hat{\mu}_{c,\delta}|^2 \leq q_d |\mu_c^c(t_c) - \hat{\mu}_{c,\delta}|^2 + 2\rho^2 M_{g_c}^2 D_x^2 + 2\rho^2 M_{g_c}^2 q_p^{2\text{ops}(k_c,t)} L_x^2 + 2\rho^2 M_{g_c}^2 D_x q_p^{\text{ops}(k_c,t)} L_x,$$

where $q_d := 3(1 - \rho\delta)^2 \in [0, 1)$, $M_{g_c} := \max_{x \in X} \|\nabla g_c(x)\|_{max}$, $D_x := \max_{x, y \in X} \|x - y\|_{max}$, $L_x = \max_j \|x^j(0) - x^*(t)\|_{max}$, t is the dual iteration count vector onboard the primal agents when sending states to dual agent c , and k_c is the time at which the first primal agent sent a state to dual agent c with $\mu(t)$ onboard.

Proof: Let $x^c(t)$ be denoted by x_t^c , and define $\hat{x}^t = \arg \min_{x \in X} L_\delta(x, \mu(t))$ and $\hat{x}_\delta = \arg \min_{x \in X} L_\delta(x, \hat{\mu}_\delta)$.

Using the non-expansiveness of Π_M , we find

$$\begin{aligned} |\mu_c^c(t_c + 1) - \hat{\mu}_{c,\delta}|^2 &= |\Pi_M[\mu_c^c(t_c) + \rho(g_c(x_t^c) - \delta\mu_c^c(t_c))] - \Pi_M[\hat{\mu}_{c,\delta} + \rho(g_c(\hat{x}_\delta) - \delta\hat{\mu}_{c,\delta})]|^2 \\ &\leq (1 - \rho\delta)^2 |\mu_c^c(t_c) - \hat{\mu}_{c,\delta}|^2 + \rho^2 |g_c(x_t^c) - g_c(\hat{x}_\delta)|^2 - 2\rho(1 - \rho\delta)(\mu_c^c(t_c) - \hat{\mu}_{c,\delta})(g_c(\hat{x}_\delta) - g_c(x_t^c)). \end{aligned}$$

Adding $g_c(\hat{x}^t) - g_c(\hat{x}^\delta)$ in the last set of parentheses gives

$$\begin{aligned} |\mu_c^c(t_c + 1) - \hat{\mu}_{c,\delta}|^2 &\leq (1 - \rho\delta)^2 |\mu_c^c(t_c) - \hat{\mu}_{c,\delta}|^2 + \rho^2 |g_c(x_t^c) - g_c(\hat{x}_\delta)|^2 \\ &\quad - 2\rho(1 - \rho\delta)(\mu_c^c(t_c) - \hat{\mu}_{c,\delta})(g_c(\hat{x}_\delta) - g_c(\hat{x}^t)) \\ &\quad - 2\rho(1 - \rho\delta)(\mu_c^c(t_c) - \hat{\mu}_{c,\delta})(g_c(\hat{x}^t) - g_c(x_t^c)). \end{aligned} \quad (4)$$

Using $0 \leq |(1 - \rho\delta)(\mu_c^c(t_c) - \hat{\mu}_{c,\delta}) + \rho(g_c(\hat{x}^t) - g_c(x_t^c))|^2$, we expand and rearrange to give

$$-2\rho(1 - \rho\delta)(\mu_c^c(t_c) - \hat{\mu}_{c,\delta})(g_c(\hat{x}^t) - g_c(x_t^c)) \leq (1 - \rho\delta)^2 |\mu_c^c(t_c) - \hat{\mu}_{c,\delta}|^2 + \rho^2 |g_c(\hat{x}^t) - g_c(x_t^c)|^2. \quad (5)$$

Similarly, we can derive

$$-2\rho(1 - \rho\delta)(\mu_c^c(t_c) - \hat{\mu}_{c,\delta})(g_c(\hat{x}_\delta) - g_c(\hat{x}^t)) \leq (1 - \rho\delta)^2 |\mu_c^c(t_c) - \hat{\mu}_{c,\delta}|^2 + \rho^2 |g_c(\hat{x}_\delta) - g_c(\hat{x}^t)|^2. \quad (6)$$

Using Equations (5) and (6) in Equation (4) gives

$$\begin{aligned} |\mu_c^c(t_c+1) - \hat{\mu}_{c,\delta}|^2 &\leq 3(1 - \rho\delta)^2 |\mu_c^c(t_c) - \hat{\mu}_{c,\delta}|^2 + \rho^2 |g_c(\hat{x}_\delta) - g_c(\hat{x}^t)|^2 \\ &\quad + \rho^2 |g_c(\hat{x}^t) - g_c(x_t^c)|^2 + \rho^2 |g_c(x_t^c) - g_c(\hat{x}_\delta)|^2. \end{aligned} \quad (7)$$

For the $\rho^2 |g_c(x_t^c) - g_c(\hat{x}_\delta)|^2$ term, we can write

$$\begin{aligned} |g_c(x_t^c) - g_c(\hat{x}_\delta)|^2 &= |g_c(x_t^c) - g_c(\hat{x}^t) + g_c(\hat{x}^t) - g_c(\hat{x}_\delta)|^2 \\ &\leq |g_c(x_t^c) - g_c(\hat{x}^t)|^2 + |g_c(\hat{x}^t) - g_c(\hat{x}_\delta)|^2 + 2|g_c(x_t^c) - g_c(\hat{x}^t)||g_c(\hat{x}^t) - g_c(\hat{x}_\delta)|, \end{aligned}$$

where substituting this into Equation (7) and grouping gives

$$\begin{aligned} |\mu_c^c(t_c + 1) - \hat{\mu}_{c,\delta}|^2 &\leq 3(1 - \rho\delta)^2 |\mu_c^c(t_c) - \hat{\mu}_{c,\delta}|^2 + 2\rho^2 |g_c(\hat{x}^t) - g_c(\hat{x}_\delta)|^2 \\ &\quad + 2\rho^2 |g_c(x_t^c) - g_c(\hat{x}^t)|^2 + 2\rho^2 |g_c(x_t^c) - g_c(\hat{x}^t)||g_c(\hat{x}^t) - g_c(\hat{x}_\delta)|. \end{aligned}$$

Using the Lipschitz property of g_c and the definition of D_x ,

$$\begin{aligned} |\mu_c^c(t_c+1) - \hat{\mu}_{c,\delta}|^2 &\leq 3(1 - \rho\delta)^2 |\mu_c^c(t_c) - \hat{\mu}_{c,\delta}|^2 + 2\rho^2 M_{g_c}^2 D_x^2 \\ &\quad + 2\rho^2 M_{g_c}^2 \|x_t^c - \hat{x}^t\|_{max}^2 + 2\rho^2 M_{g_c}^2 D_x \|x_t^c - \hat{x}^t\|_{max}. \end{aligned} \quad (8)$$

Previously, we established primal convergence for a fixed $\mu(t)$. This allows us to write $\|x_t^c - \hat{x}^t\|_{max} \leq q_p^{\text{ops}(k_c, t)} L_x$.

Substituting this into Equation (8) completes the proof. \blacksquare

Remark 1: The term $2\rho^2 M_{g_c}^2 D_x^2$ in Theorem 3 is termed the ‘‘asynchrony penalty,’’ because it is an offset from reaching a solution. It is due to asynchronously computing dual variables and is absent when dual updates are centralized [12], [15]. Further bounding it will be the subject of future research.

We next present a simplified dual convergence rate.

Theorem 4: Let all conditions of Theorem 3 hold. In Algorithm 2, convergence for dual agent c obeys

$$|\mu_c^c(t_c + 1) - \hat{\mu}_{c,\delta}|^2 \leq q_d^{t_c+1} |\mu_c^c(0) - \hat{\mu}_{c,\delta}|^2 + \frac{1 - q_d^{t_c+2}}{1 - q_d} p,$$

where $p = \max_t 2\rho^2 M_{g_c}^2 D_x^2 + 2\rho^2 M_{g_c}^2 q_p^{2\text{ops}(k_c,t)} L_x^2 + 2\rho^2 M_{g_c}^2 D_x q_p^{\text{ops}(k_c,t)} L_x$.

Proof: First, define

$$p_i := 2\rho^2 M_{g_c}^2 D_x^2 + 2\rho^2 M_{g_c}^2 q_p^{2\text{ops}(k_c,i)} L_x^2 + 2\rho^2 M_{g_c}^2 D_x q_p^{\text{ops}(k_c,i)}.$$

Then, recursively applying Theorem 3 gives

$$|\mu_c^c(t_c + 1) - \hat{\mu}_{c,\delta}|^2 \leq q_d^{t_c+1} |\mu_c^c(0) - \hat{\mu}_{c,\delta}|^2 + \sum_{j=0}^{t_c+1} q_d^{t_c+1-j} p_j. \quad (9)$$

By definition, $q_d \in [0, 1)$, so, using $p_i \leq p$, summing the geometric series in Equation (9) completes the proof. ■

Theorems 3 and 4 can be used to give an overall primal convergence rate for $x^i(k; t)$ converging to \hat{x}_δ .

Theorem 5: Let Assumptions 1-4 hold. Then there exist constants $K_1, K_2 > 0$ such that, for all t ,

$$\|x^i(k; t) - \hat{x}_\delta\|_2 \leq K_1 q_p^{\text{ops}(k,t)} + K_2 \|\mu(t) - \hat{\mu}_\delta\|_2.$$

Proof: Plug Theorem 3 into Theorem 2 of [12]. ■

V. NUMERICAL EXAMPLE

We consider an example with $n = 10$ primal agents (each updating a scalar variable) whose objective function is

$$h(x) = \sum_{i=1}^n x_i^4 + \frac{1}{20} \sum_{i=1}^n \sum_{\substack{j=1 \\ j \neq i}}^n (x_i - x_j)^2.$$

Set $\delta = 0.001$. For $b = (-2, 4, -10, 5, 1, 8)^T$ and

$$A = \begin{pmatrix} -1 & 0 & -3 & 0 & 0 & 4 & 0 & 0 & 10 & 0 \\ 0 & 1 & 5 & 1 & 1 & 0 & 0 & 2 & 0 & 5 \\ 0 & 0 & 1 & 1 & -5 & 1 & 4 & 0 & 0 & 0 \\ 0 & 0 & -2 & 0 & 0 & 8 & 1 & 1 & -3 & 1 \\ 0 & 0 & 0 & 0 & -3 & 0 & 1 & 1 & 1 & 0 \\ 0 & 4 & 0 & 0 & 0 & 0 & 0 & 2 & 1 & -4 \end{pmatrix},$$

there are $m = 6$ dual agents each responsible for a scalar dual variable that encodes a constraint in $g(x) = Ax - b$.

We also require that each $x_i \in X_i = [1, 10]$. Then

$$L_\delta(x, \mu) = \sum_{i=1}^n x_i^4 + \frac{1}{20} \sum_{i=1}^n \sum_{\substack{j=1 \\ j \neq i}}^n (x_i - x_j)^2 + \mu^T (Ax - b) - \frac{\delta}{2} \|\mu\|^2.$$

We find that H is β -diagonally dominant with $\beta = 12$.

We use this simulation example to explore fundamental relationships between the diagonal dominance parameter β and the communication rate between agents. Here, we use random communications and vary the probability that

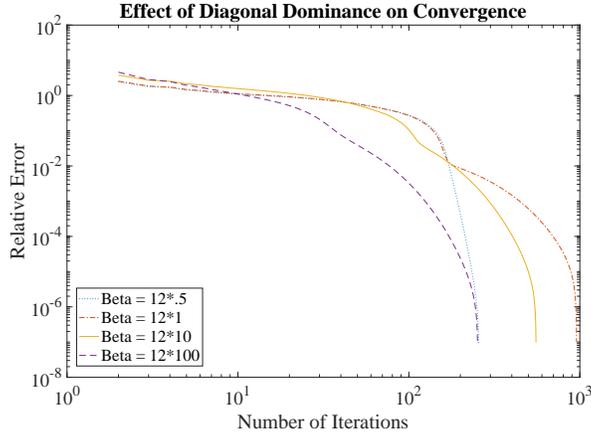


Fig. 1. Effect of diagonal dominance on convergence. Larger values of β tend to produce faster convergence. However, sufficiently small values for β may also produce the same result by promoting constraint satisfaction earlier.

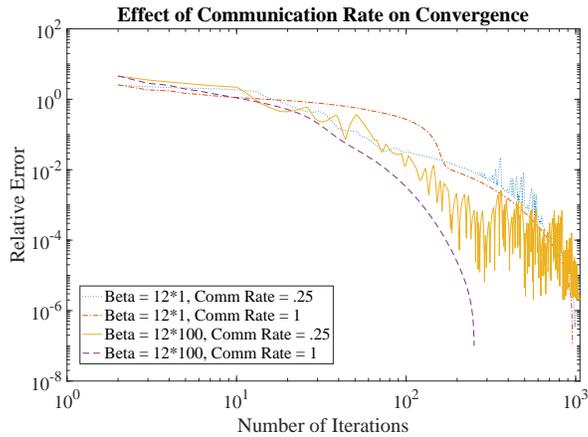


Fig. 2. Effect of communication rate on convergence. Less frequent communication leads to slower convergence, larger errors, and oscillation in the optimum's proximity. The benefits of strong diagonal dominance can also be outweighed by poor communication.

agents i and j communicate at a particular timestep. We begin by varying β over $\beta \in 12 \cdot \{.9, 1, 10, 100\}$, which we do by scaling h while holding other terms constant. Figure 1 plots the iteration number k versus the value $\frac{\|x_i^i(k;t) - \hat{x}_\delta\|}{\|\hat{x}_\delta\|}$ (the relative error) for these values of β and a communication rate of 1 (the agents communicate every update with one another).

As predicted by Theorem 2, a larger β correlates with faster convergence in general. Figure 1 also reveals, however, that if β is sufficiently small, then the convergence rate is eventually similar to that for a large β . This is explained by how β weights h versus g : if β is large, then h will be minimized quickly and satisfaction of g will only be attained afterwards, which prolongs convergence. Conversely, for smaller β , minimizing h and satisfying g are weighted more equally, which promotes convergence of both equally.

Varying the communication rate has a significant impact on the number of iterations required and the behavior

of error in agents' updates. Communication rate can even outweigh the benefits of a large β , shown in Figure 2. This reveals that faster convergence can be achieved by both improving communication or increasing the diagonal dominance of the problem. However, if communication is poor, increasing diagonal dominance may not improve results by much, which suggests that even favorable problem structure does not eliminate the impact of asynchrony. Oscillations in the proximity of the optimum are also amplified by large values of β .

VI. CONCLUSION

After exploring a counterexample to asynchronous dual communications, Algorithm 2 presents an asynchronous primal-dual approach that is asynchronous in primal updates and communications and asynchronous in distributed dual updates. A numerical example illustrates the effect diagonal dominance has with other parameters upon convergence. Future work will apply the algorithm to large-scale machine learning problems and explore reducing the asynchrony penalty.

APPENDIX

Proof of Theorem 1: Consider the quadratic program

$$\begin{aligned} & \text{minimize } \frac{1}{2}x^T Qx + r^T x \\ & \text{subject to } Ax \leq b, \quad x \in X, \end{aligned}$$

where $\mathbb{R}^{n \times n} \ni Q = Q^T \succ 0$, $r \in \mathbb{R}^n$, $A \in \mathbb{R}^{m \times n}$, and $b \in \mathbb{R}^m$. We take X sufficiently large in a sense to be described below.

Because Q is symmetric and positive definite, its eigenvectors can be orthonormalized, and we denote these eigenvectors by v_1, \dots, v_n . To construct an example, suppose that $A \in \mathbb{R}^{m \times n}$ has row i equal to the i^{th} normalized eigenvector of Q . To compute $\hat{x}_1 := \arg \min_{x \in X} L_\delta(x, \mu^1)$ we set $\nabla_x L_\delta(\hat{x}_1, \mu^1) = 0$. Expanding and solving, we find $\hat{x}_1 = -Q^{-1}A^T \mu^1$, where we assume that X is large enough to contain this point. Similarly, we find $\hat{x}_2 = -Q^{-1}A^T \mu^2$, where we also assume $\hat{x}_2 \in X$. Then

$$\|\hat{x}_1 - \hat{x}_2\|_2 = \|Q^{-1}A^T(\mu^2 - \mu^1)\|_2 \geq \sigma_{\min}(Q^{-1}A^T)\|\mu^1 - \mu^2\|_2,$$

where $\sigma_{\min}(\cdot)$ is the minimum singular value. Expanding,

$$\sigma_i^2(Q^{-1}A^T) = \lambda_i(AQ^{-T}Q^{-1}A^T) = \lambda_i(AQ^{-2}A^T),$$

and $AQ^{-2}A^T = \text{diag}\left(\frac{1}{\lambda_1(Q)^2}, \dots, \frac{1}{\lambda_n(Q)^2}\right)$, which follows from the fact that we have orthonormalized Q 's eigenvectors. Then $\sigma_{\min}(Q^{-1}A^T) = \frac{1}{\lambda_{\max}(Q)}$. Using this above, we find

$$\|\hat{x}_1 - \hat{x}_2\|_2 \geq \frac{1}{\lambda_{\max}(Q)}\|\mu^1 - \mu^2\|_2.$$

To enforce $\|\hat{x}_1 - \hat{x}_2\|_2 > L$, we ensure that $\frac{\epsilon}{\lambda_{\max}(Q)} > L$, which is attained for any matrix satisfying $\lambda_{\max}(Q) < \frac{\epsilon}{L}$.

■

REFERENCES

- [1] Kenneth Joseph Arrow, Hirofumi Azawa, Leonid Hurwicz, and Hirofumi Uzawa. *Studies in linear and non-linear programming*, volume 2. Stanford University Press, 1958.
- [2] Dimitri P. Bertsekas. Distributed asynchronous computation of fixed points. *Mathematical Programming*, 27(1):107–120, Sep 1983.
- [3] Dimitri P Bertsekas and Athena Scientific. *Convex optimization algorithms*. Athena Scientific Belmont, 2015.
- [4] Dimitri P. Bertsekas and John N. Tsitsiklis. *Parallel and Distributed Computation: Numerical Methods*. Prentice-Hall, Inc., USA, 1989.
- [5] Dimitri P. Bertsekas and John N. Tsitsiklis. Some aspects of the parallel and distributed iterative algorithms—a survey. *Automatica*, 27(1):3–21, January 1991.
- [6] Stephen Boyd, Stephen P Boyd, and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- [7] John C Duchi, Alekh Agarwal, and Martin J Wainwright. Dual averaging for distributed optimization: Convergence analysis and network scaling. *IEEE Transactions on Automatic control*, 57(3):592–606, 2011.
- [8] Francisco Facchinei and Jong-Shi Pang. *Finite-dimensional variational inequalities and complementarity problems*. Springer Science & Business Media, 2007.
- [9] Derek Greene and Pádraig Cunningham. Practical solutions to the problem of diagonal dominance in kernel document clustering. In *Proceedings of the 23rd International Conference on Machine Learning*, page 377384, 2006.
- [10] M. T. Hale and M. Egerstedt. Cloud-based optimization: A quasi-decentralized approach to multi-agent coordination. In *53rd IEEE Conference on Decision and Control*, pages 6635–6640, 2014.
- [11] M. T. Hale, A. Nedi, and M. Egerstedt. Cloud-based centralized/decentralized multi-agent optimization with communication delays. In *2015 54th IEEE Conference on Decision and Control (CDC)*, pages 700–705, Dec 2015.
- [12] M. T. Hale, A. Nedi, and M. Egerstedt. Asynchronous multiagent primal-dual optimization. *IEEE Transactions on Automatic Control*, 62(9):4421–4435, Sep. 2017.
- [13] S. Hochhaus and M. T. Hale. Asynchronous distributed optimization with heterogeneous regularizations and normalizations. In *2018 IEEE Conference on Decision and Control (CDC)*, pages 4232–4237, 2018.
- [14] Martin Jaggi, Virginia Smith, Martin Takáč, Jonathan Terhorst, Sanjay Krishnan, Thomas Hofmann, and Michael I Jordan. Communication-efficient distributed dual coordinate ascent. In *Advances in neural information processing systems*, pages 3068–3076, 2014.
- [15] Jayash Koshal, Angelia Nedić, and Uday V Shanbhag. Multiuser optimization: Distributed algorithms and error analysis. *SIAM Journal on Optimization*, 21(3):1046–1081, 2011.
- [16] A. Nedic, A. Ozdaglar, and P. A. Parrilo. Constrained consensus and optimization in multi-agent networks. *IEEE Transactions on Automatic Control*, 55(4):922–938, 2010.
- [17] Angelia Nedic and Asuman Ozdaglar. Distributed subgradient methods for multi-agent optimization. *IEEE Transactions on Automatic Control*, 54(1):48–61, 2009.
- [18] A. Nedi and A. Olshevsky. Distributed optimization over time-varying directed graphs. *IEEE Transactions on Automatic Control*, 60(3):601–615, 2015.
- [19] Ivano Notarnicola and Giuseppe Notarstefano. Asynchronous distributed optimization via randomized dual proximal gradient. *IEEE Transactions on Automatic Control*, 62(5):2095–2106, 2016.
- [20] Shai Shalev-Shwartz. Online learning and online convex optimization. *Foundations and Trends in Machine Learning*, 4(2):107–194, 2012.
- [21] Suvrit Sra, Sebastian Nowozin, and Stephen J Wright. *Optimization for machine learning*. Mit Press, 2012.
- [22] K. I. Tsianos, S. Lawlor, and M. G. Rabbat. Push-sum distributed dual averaging for convex optimization. In *2012 IEEE 51st IEEE Conference on Decision and Control (CDC)*, pages 5453–5458, 2012.
- [23] K. I. Tsianos and M. G. Rabbat. Distributed dual averaging for convex optimization under communication delays. In *2012 American Control Conference (ACC)*, pages 1067–1072, 2012.
- [24] John Tsitsiklis, Dimitri Bertsekas, and Michael Athans. Distributed asynchronous deterministic and stochastic gradient optimization algorithms. *IEEE transactions on automatic control*, 31(9):803–812, 1986.
- [25] M. Ubl and M. T. Hale. Totally asynchronous distributed quadratic programming with independent stepsizes and regularizations. In *58th Conference on Decision and Control (CDC)*, pages 7423–7428, 2019.
- [26] Minghui Zhu and Sonia Martínez. On distributed convex optimization under inequality and equality constraints. *IEEE Transactions on Automatic Control*, 57(1):151–164, 2011.