

A primal-dual method for conic constrained distributed optimization problems

Authored by:

Necdet Serhat Aybat
Erfan Yazdandoost Hamedani

Abstract

We consider cooperative multi-agent consensus optimization problems over an undirected network of agents, where only those agents connected by an edge can directly communicate. The objective is to minimize the sum of agent-specific composite convex functions over agent-specific private conic constraint sets; hence, the optimal consensus decision should lie in the intersection of these private sets. We provide convergence rates in sub-optimality, infeasibility and consensus violation; examine the effect of underlying network topology on the convergence rates of the proposed decentralized algorithms; and show how to extend these methods to handle time-varying communication networks.

1 Paper Body

Let $G = (N, E)$ denote a connected undirected graph of N computing nodes, where $N = \{1, \dots, N\}$ and $E \subseteq N \times N$ denotes the set of edges without loss of generality assume that $(i, j) \in E$ implies $i \leq j$. Suppose nodes i and j can exchange information only if $(i, j) \in E$, and each node $i \in N$ has a private (local) cost function $\varphi_i : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ such that $\varphi_i(x) = \varphi_i(x) + f_i(x)$,

(1)

where $\varphi_i : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is a possibly non-smooth convex function, and $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$ is a smooth convex function. We assume that f_i is differentiable on an open set containing $\text{dom } \varphi_i$ with a Lipschitz continuous gradient ∇f_i , of which Lipschitz constant is L_i ; and the prox map of φ_i ,

$$\text{prox}_{\varphi_i}(x) = \arg\min_{y \in \mathbb{R}^n} \varphi_i(y) + \frac{1}{2} \|x - y\|^2,$$

(2)

$$\|x\|_2^2,$$

(2)

is efficiently computable for $i \in N$, where $\| \cdot \|$ denotes the Euclidean norm. Let $N_i = \{j \in N : (i, j) \in E \text{ or } (j, i) \in E\}$ denote the set of neighboring nodes of $i \in N$, and $d_i = |N_i|$ is the degree of node $i \in N$. Consider the following minimization problem: \min

$$\begin{aligned}
& x \in \mathbb{R}^n \\
& X \\
& i \in N \\
& f_i(x) \\
& \text{s.t.} \\
& A_i x \leq b_i \in K_i, \\
& i \in N, \\
& (3)
\end{aligned}$$

where $A_i \in \mathbb{R}^{m_i \times n}$, $b_i \in \mathbb{R}^{m_i}$ and $K_i \subset \mathbb{R}^{m_i}$ is a closed, convex cone. Suppose that projections onto K_i can be computed efficiently, while the projection onto the preimage $A_i^{-1}(K_i + b_i)$ is assumed to be impractical, e.g., when K_i is the positive semidefinite cone, projection to preimage requires solving an SDP. Our objective is to solve (3) in a decentralized fashion using the computing nodes N and exchanging information only along the edges E . In Section 2 and Section 3, we consider (3) when the topology of the connectivity graph is static and time-varying, respectively. This computational setting, i.e., decentralized consensus optimization, appears as a generic model for various applications in signal processing, e.g., [1, 2], machine learning, e.g., [3, 4, 5] and statistical inference, e.g., [6]. Clearly, (3) can also be solved in a ‘centralized’ fashion by communicating all the private functions f_i to a central node, and solving the overall problem at this 30th Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain.

node. However, such an approach can be very expensive both from communication and computation perspectives when compared to the distributed algorithms which are far more scalable to increasing problem data and network sizes. In particular, suppose $(A_i, b_i) \in \mathbb{R}^{m_i \times (n+1)}$ and $2 f_i(x) = \|x\|_1 + \|A_i x - b_i\|_1$ for some given $\lambda \geq 0$ for $i \in N$ such that $m \leq n$ and $N \geq 1$. Hence, (3) is a very large scale LASSO problem with distributed data. To solve (3) in a centralized fashion, the data $\{(A_i, b_i) : i \in N\}$ needs to be communicated to the central node. This can be prohibitively expensive, and may also violate privacy constraints if in case some node i does not want to reveal the details of its private data. Furthermore, it requires that the central node has large enough memory to be able to accommodate all the data. On the other hand, at the expense of slower convergence, one can completely do away with a central node, and seek for consensus among all the nodes on an optimal decision using ‘local’ decisions communicated by the neighboring nodes. From computational perspective, for certain cases, computing partial gradients locally can be more computationally efficient when compared to computing the entire gradient at a central node. With these considerations in mind, we propose decentralized algorithms that can compute solutions to (3) using only local computations without explicitly requiring the nodes to communicate the functions $\{f_i : i \in N\}$; thereby, circumventing all privacy, communication and memory issues. Examples of constrained machine learning problems that fit into our framework include multiple kernel learning [7], and primal linear support vector machine (SVM) problems. In the numerical section we implemented the proposed algorithms on the primal SVM problem. 1.1

Previous Work

There has been active research [8, 9, 10, 11, 12] on solving convex-concave saddle point problems $\min_x \max_y L(x, y)$. In [9] primal-dual proximal algorithms are proposed for convex-concave problems with known saddle-point structure $\min_x \max_y L_s(x, y)$, $\varphi(x) + h^T x, y \leq h(y)$, where φ and h are convex functions, and T is a linear map. These algorithms converge with rate $O(1/k)$ for the primal-dual gap, and they can be modified to yield a convergence rate of $O(1/k^2)$ when either φ or h is strongly convex, and $O(1/ek)$ linear rate, when both φ and h are strongly convex. More recently, in [11] Chambolle and Pock extend their previous work in [9], using simpler proofs, to handle composite convex primal functions, i.e., sum of smooth and (possibly) nonsmooth functions, and to deal with proximity operators based on Bregman distance functions. Consider $\min_{x \in \mathbb{R}^n} \sum_{i=1}^N \varphi_i(x) : x \in \bigcap_{i=1}^N X_i$ over $G = (N, E)$. Although the unconstrained consensus optimization, i.e., $X_i = \mathbb{R}^n$, is well studied (see [13, 14] and the references therein, the constrained case is still an immature, and recently developing area of active research [13, 14, 15, 16, 17, 18, 19]. Other than few exceptions, e.g., [15, 16, 17], the methods in literature require that each node compute a projection on the privately known set X_i in addition to consensus and (sub)gradient steps, e.g., [18, 19]. Moreover, among those few exceptions that do not use projections onto X_i when φ_i is not easy to compute, only [15, 16] can handle agent-specific constraints without assuming global knowledge of the constraints by all agents. However, no rate results in terms of suboptimality, local infeasibility, and consensus violation exist for the primaldual distributed methods in [15, 16] when implemented for the agent-specific conic constraint sets $X_i = \{x : A_i x \leq b_i, x \in K_i\}$ studied in this paper. In [15], a consensus-based distributed primaldual perturbation (PDP) algorithm using a square summable but not summable step-size sequence is proposed. The objective is to minimize a composition of a global network function (smooth) with the summation of local objective functions (smooth), subject to local compact sets and inequality constraints on the summation of agent specific constrained functions. They showed that the local primal-dual iterate sequence converges to a global optimal primal-dual solution; however, no rate result was provided. The proposed PDP method can also handle non-smooth constraints with similar convergence guarantees. Finally, while we were preparing this paper, we became aware of a very recent work [16] related to ours. The authors proposed a distributed algorithm on time-varying communication network for solving saddle-point problems subject to consensus constraints. The algorithm can also be applied to solve consensus optimization problems with inequality constraints that can be written as summation of local convex functions of local and global variables. Under some assumptions, it is shown that using a carefully selected decreasing step-size sequence, the ergodic average of primal-dual sequence converges with $O(1/k)$ rate in terms of saddle-point evaluation error; however, when applied to constrained optimization problems, no rate in terms of either suboptimality or infeasibility is provided.

2

Contribution. We propose primal-dual algorithms for distributed optimiza-

First, we define the notation used throughout the paper. Next, in Theorem 1.1, we discuss a special case of (4), which will help us prove the main results of this paper, and also allow us to develop decentralized algorithms for the consensus optimization problem in (3). The proposed algorithms in this paper can distribute the computation over the nodes such that each node's computation is based on the local topology of G and the private information only available to that node.

Notation. Throughout the paper, $\| \cdot \|$ denotes the Euclidean norm. Given a convex set S , let $\mathcal{S}(\cdot)$ denote its support function, i.e., $\mathcal{S}(\cdot) = \sup_{w \in S} \langle \cdot, w \rangle$, let $\mathcal{I}_S(\cdot)$ denote the indicator function of S , i.e., $\mathcal{I}_S(w) = 0$ for $w \in S$ and equal to $+\infty$ otherwise, and let $\mathcal{P}_S(w) = \arg\min\{\|v - w\| : v \in S\}$ denote the projection onto S . For a closed convex set S , we define the distance function as $d_S(w) = \mathcal{K}_S(w) = \|w\|_K$. Given a convex cone $K \subseteq \mathbb{R}^m$, let K° denote its dual cone, i.e., $K^\circ = \{v \in \mathbb{R}^m : \langle v, w \rangle \geq 0, \forall w \in K\}$, and $K^\circ = K^\circ$ denote the polar cone of K . Note that for a given cone $K \subseteq \mathbb{R}^m$, $\mathcal{K}(\cdot) = 0$ for $\cdot \in K$ and equal to $+\infty$ if $\cdot \notin K$, i.e., $\mathcal{K}(\cdot) = \mathcal{I}_K(\cdot)$ for all $\cdot \in \mathbb{R}^m$. Cone K is called proper if it is closed, convex, pointed, and it has a nonempty interior. Given a convex function $g : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$, its convex conjugate is defined as $g^*(w) = \sup_{x \in \mathbb{R}^n} \langle w, x \rangle - g(x)$. \otimes denotes the Kronecker product, and I_n is the $n \times n$ identity matrix. **Definition 1.** Let $X, Y \subseteq \mathbb{R}^n$ and $X = [x_i]_{i=1}^n$; $Y = [y_i]_{i=1}^n$, $x_i \in \mathbb{R}^{m_i}$, $y_i \in \mathbb{R}^{n_i}$, $y = P[x, z]$ and $z = [z_i]_{i=1}^n \in \mathbb{R}^{m_i}$, where m_i, n_i and n denotes the Cartesian product. Given parameters $\alpha_i \geq 0, \beta_i \geq 0$ for $i = 1, \dots, n$, let $D^x = \text{diag}(\alpha_i I_{m_i})$, $D^y = \text{diag}(\beta_i I_{n_i})$, and $D = \text{diag}([\alpha_i I_{m_i} \beta_i I_{n_i}]_{i=1}^n)$. Defining $\|x\|_x = \sqrt{x^T D^x x}$ and $\|y\|_y = \sqrt{y^T D^y y}$, $\|z\|_z = \sqrt{z^T D^z z}$ leads to the following Bregman distance functions: $D_x(x, x')$

$\|z\|_z = \sqrt{z^T D^z z}$, where the Q -norm is defined as $\|z\|_Q = \sqrt{z^T Q z}$ for $Q \succeq 0$.

D^x
2
 Q
 D^y
3

Theorem 1.1. Let X, Y , and Bregman functions D_x, D_y be defined as in Definition 1. Suppose $P = [P_i]_{i=1}^n$, $P_i(x_i) = h_0(\cdot) + \sum_{i=1}^n h_i(\cdot)$, where $\{h_i\}_{i=1}^n$ are composite convex functions defined as in (1), and $\{h_i\}_{i=1}^n$ are closed convex with simple prox-maps. Given $A_0 \in \mathbb{R}^{n \times n}$ and $\{A_i\}_{i=1}^n$ such that $A_i \in \mathbb{R}^{m_i \times n_i}$, let $T = [A \ A_0]$, where $A = \text{diag}([A_i]_{i=1}^n) \in \mathbb{R}^{(m+n) \times (m+n)}$ is a block-diagonal matrix. Given the initial point (x_0, y_0) , the PDA iterate sequence according to $\{x_k, y_k\}_{k=1}^K$, generated by (5a) and (5b) when $\alpha_x = \alpha_y = 1$, satisfies (6) for all $K \geq 1$.

if $Q(A, A_0) \succeq 0$

A_0

$A^T D^x 0$

$A^T 0 \succeq 0$, $\text{diag}([(\alpha_i - \beta_i) I_{n_i}]_{i=1}^n)$. Moreover, if $\alpha \geq 0$, where $D = \alpha D^x$

If a saddle point exists for (4), and $Q(A, A_0) \neq 0$, then $\{x_k, y_k\}_{k=1}^\infty$ converges to a saddle point of (4); hence, $\{x, y\}_{k=1}^\infty$ converges to the same point.

Although the proof of Theorem 1.1 follows from the lines of [11], we provide the proof in the appendix for the sake of completeness as it will be used repeatedly to derive our results. Next we discuss how (5) can be implemented to compute an ϵ -optimal solution to (3) in a distributed way using only $O(1/\epsilon)$ communications over the communication graph G while respecting nodespecific privacy requirements. Later, in Section 3, we consider the scenario where the topology of the connectivity graph is time-varying, and propose a distributed algorithm that requires $O(1/\epsilon + p)$ communications for any $p \geq 1$. Finally, in Section 4 we test the proposed algorithms for solving the primal SVM problem in a decentralized manner. These results are shown under Assumption 1.1. Assumption 1.1. The duality gap for (3) is zero, and a primal-dual solution to (3) exists. A sufficient condition for this is the existence of a Slater point, i.e., there exists $x \in \text{relint}(\text{dom } \varphi) \cap \bigcap_{i=1}^N \text{int}(K_i)$ for $i \in N$, where $\text{dom } \varphi = \bigcap_{i=1}^N \text{dom } \varphi_i$ such that $A_i x$

2

Static Network Topology

Let $x_i \in \mathbb{R}^n$ denote the local decision vector of node $i \in N$. By taking advantage of the fact that G is connected, we can reformulate (3) as the following distributed consensus optimization problem: \min

$x_i \in \mathbb{R}^n, i \in N$

(

X

$i \in N$

$\varphi_i(x_i) - x_i = x_j : \varphi_{ij}, (i, j) \in E,$

$A_i x_i \in b_i \cap K_i : \varphi_i, i \in N$

)

,

(7)

where $\varphi_{ij} \in \mathbb{R}^n$ and $\varphi_i \in \mathbb{R}^{m_i}$ are the corresponding dual variables. Let $x = [x_i]_{i \in N} \in \mathbb{R}^{n \times N}$. The consensus constraints $x_i = x_j$ for $(i, j) \in E$ can be formulated as $Mx = 0$, where $M \in \mathbb{R}^{n \times E \times n \times N}$ is a block matrix such that $M = H \otimes I_n$ where H is the oriented edge-node incidence matrix, i.e., the entry $H(i, j), l$, corresponding to edge $(i, j) \in E$ and $l \in N$, is equal to 1 if $l = i$, -1 if $l = j$, and 0 otherwise. Note that $M^T M = H^T H \otimes I_n = \mathcal{L} \otimes I_n$, where $\mathcal{L} \in \mathbb{R}^{n \times n}$ denotes the graph Laplacian of G , i.e., $\mathcal{L}_{ii} = d_i$, $\mathcal{L}_{ij} = -1$ if $(i, j) \in E$ or $(j, i) \in E$, and equal to 0 otherwise. For any closed convex set S , we have $\mathcal{L}S \subseteq \mathcal{L}S$; therefore, using the fact that $\mathcal{L}K = \bigcap_{i \in N} K_i$ for $i \in N$, one can obtain the following saddle point problem corresponding to (7),

$\min \max L(x, y), x$

y

X

$i \in N$

$\varphi_i(x_i) + h^T \varphi_i, A_i x_i \in b_i \cap K_i : \varphi_i \in h^T, Mx = 0,$

where $y = [y_1 \dots y_m]^T$ for $y_i = [y_{ij}]_{j \in \mathcal{N}_i} \in \mathbb{R}^{n_i - 1}$, $y_i = [y_{ij}]_{j \in \mathcal{N}_i} \in \mathbb{R}^{m_i}$, and $m = \sum_{i \in \mathcal{N}} m_i$.

$$\begin{aligned} &P \\ & \mathcal{N} \\ & (8) \\ & m_i. \end{aligned}$$

Next, we study the distributed implementation of PDA in (5a)-(5b) to solve (8). Let $f(x) = P^T P^T \sum_{i \in \mathcal{N}} f_i(x_i)$, and $h(y) = \sum_{i \in \mathcal{N}} K_i(y_i) + h(y)$, $y_i \in \mathbb{R}^{m_i}$. Define the block-diagonal matrix $A = \text{diag}([A_i]_{i \in \mathcal{N}}) \in \mathbb{R}^{m \times m}$ and $T = [A^T M^T]^T$. Therefore, given the initial iterates x_0, y_0 and parameters $\alpha \in (0, 1)$, $\beta_i \in (0, 1)$ for $i \in \mathcal{N}$, choosing D_x and D_y as defined in Definition 1, and setting $\alpha_x = \alpha_y = 1$, PDA iterations in (5a)-(5b) take the following form: $h(x_{k+1}) = \arg\min_{x \in \mathcal{X}} h(x) + \sum_{i \in \mathcal{N}} \lambda_i (x_i - x_{k,i})^T$

$$\begin{aligned} &X \\ & \mathcal{N} \\ & f_i(x_i) + h(f(x_k)), x_{k,i} + h(A_i x_k - b_i), \lambda_{k,i} + \\ & \lambda_{k,i} - \lambda_{k,i}^2, \lambda_{k,i}^2 \\ & \lambda_{k,i} - \lambda_{k,i}^2, i \in \mathcal{N} \end{aligned} \quad \arg\min_{x \in \mathcal{X}} h(x) + \sum_{i \in \mathcal{N}} \lambda_i (x_i - x_{k,i})^T$$

$$(9a)$$

$$\lambda_{k+1} = \arg\min_{\lambda \in \mathcal{K}} \sum_{i \in \mathcal{N}} \lambda_i (x_{k+1,i} - x_{k,i})^T + h(x_{k+1})$$

$$(9b)$$

$$\lambda_{k+1}$$

$$(9c)$$

$$\lambda_i$$

$$4$$

Since K_i is a cone, $\text{prox}_{K_i}(\cdot) = P_{K_i}(\cdot)$; hence, λ_{k+1} can be written in closed form as

$$\begin{aligned} &\lambda_{k+1} \\ &= P_{K_i}(\lambda_{k,i} + \lambda_i A_i (x_{k+1,i} - x_{k,i}) - b_i), i \\ & \in \mathcal{N}. \end{aligned}$$

Using recursion in (9c), we can write λ_{k+1} as a partial summation of primal iterates $\{x_k\}_{k=0}^\infty$, i.e., $\lambda_{k+1} = \lambda_0 + \sum_{i=0}^k \lambda_i M (x_{i+1} - x_i)$. Let $\lambda_0 = \lambda_0 x_0$, $s_0 = x_0$, and $s_k, x_k + \lambda_k = x_k$ for $k \geq 1$; hence, $\lambda_k = \lambda_0 s_k$. Using the fact that $M^T M = I_n$, we obtain $h(x_k) = h(x_0) + \sum_{i=0}^{k-1} \lambda_i (x_{i+1} - x_i)^T$

$$P$$

$$\mathcal{N} \times \mathcal{N}$$

$$P$$

$$k \in \mathcal{N}$$

$$\lambda_{sk} \lambda_i$$

Thus, PDA iterations given in (9) for the static graph G can be computed in a decentralized way, via the node-specific computations as in Algorithm DPDA-S displayed in Fig. 1 below. Algorithm DPDA-S ($x_0, y_0, \lambda, \{\lambda_i, \lambda_i\}_{i \in \mathcal{N}}$)

Initialization: $s_0 = x_0, i \in \mathcal{N}$ Step k : ($k \geq 0$)

$$\begin{aligned}
& \mathbf{y} = [\mathbf{y}_1 \ \mathbf{y}_2 \ \dots \ \mathbf{y}_N]^\top \text{ such that } \mathbf{y}_i \in \mathbb{R}^{n_i}, \mathbf{y}_i = [\mathbf{y}_{i1} \ \mathbf{y}_{i2} \ \dots \ \mathbf{y}_{in_i}]^\top \\
& \mathbf{P} \\
& \mathbf{y}_i \\
& \mathbf{m}_i. \\
& (11)
\end{aligned}$$

Next, we consider the implementation of PDA in (5) to solve (11). Let $\mathbf{f}(\mathbf{x})$, $\mathbf{f}_i(\mathbf{x}_i)$, and $\mathbf{P}_i(\mathbf{y}_i)$ be defined as $\mathbf{f}(\mathbf{x}) = \mathbf{f}_i(\mathbf{x}_i) + \mathbf{h}_i(\mathbf{y}_i)$. Define the block-diagonal matrix $\mathbf{A} = \text{diag}([\mathbf{A}_i]_{i=1}^N) \in \mathbb{R}^{m \times n}$ and $\mathbf{T} = [\mathbf{A}^\top \ \mathbf{I}_n]^\top \in \mathbb{R}^{(m+n) \times n}$. Therefore, given the initial iterates $\mathbf{x}_0, \mathbf{y}_0$ and parameters α, β , choosing \mathbf{D}_x and \mathbf{D}_y as defined in Definition 1, and setting $\mathbf{x} = \mathbf{y} = \mathbf{0}$, PDA iterations given in (5) take the following form: Starting from $\mathbf{y}_0 = \mathbf{0}$, compute for $i = 1, 2, \dots$

$$\begin{aligned}
& \mathbf{x}_{k+1} = \argmin_{\mathbf{x}_i} \mathbf{f}_i(\mathbf{x}_i) + \mathbf{h}_i(\mathbf{x}_i) \\
& \mathbf{y}_{k+1} = \argmin_{\mathbf{y}_i} \mathbf{f}_i(\mathbf{x}_i) + \mathbf{h}_i(\mathbf{y}_i) \\
& \mathbf{x}_{k+1} = \argmin_{\mathbf{x}_i} \mathbf{f}_i(\mathbf{x}_i) + \mathbf{h}_i(\mathbf{y}_{k+1}) \\
& \mathbf{y}_{k+1} = \argmin_{\mathbf{y}_i} \mathbf{f}_i(\mathbf{x}_{k+1}) + \mathbf{h}_i(\mathbf{y}_i) \\
& (12a) \quad (12b) \\
& \mathbf{y}_{k+1} = \mathbf{y}_{k+1}. \\
& (12c)
\end{aligned}$$

Using extended Moreau decomposition for proximal operators, \mathbf{y}_{k+1} can be written as $\mathbf{y}_{k+1} = \mathbf{y}_k + \mathbf{p}(\mathbf{y}_k + \mathbf{A}^\top(\mathbf{x}_{k+1} - \mathbf{x}_k))$ where $\mathbf{p}(\mathbf{y}) = \argmin_{\mathbf{y}_i} \mathbf{f}_i(\mathbf{x}_i) + \mathbf{h}_i(\mathbf{y}_i)$.

$$\begin{aligned}
& \mathbf{y}_{k+1} = \mathbf{y}_k + \mathbf{p}(\mathbf{y}_k + \mathbf{A}^\top(\mathbf{x}_{k+1} - \mathbf{x}_k)) \\
& \mathbf{y}_{k+1} = \argmin_{\mathbf{y}_i} \mathbf{f}_i(\mathbf{x}_i) + \mathbf{h}_i(\mathbf{y}_i)
\end{aligned}$$

Let $\mathbf{1} \in \mathbb{R}^n$ be the vector all ones, $\mathbf{B}_0 = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\|_2 \leq B\}$. Note $\mathbf{P}_{\mathbf{B}_0}(\mathbf{x}) = \mathbf{x} \min\{1, \|\mathbf{x}\|_2/B\}$.

For any $\mathbf{x} = [\mathbf{x}_i]_{i=1}^N \in \mathbb{R}^n$, $\mathbf{P}_{\mathbf{B}_0}(\mathbf{x})$ can be computed as $\mathbf{P}_{\mathbf{B}_0}(\mathbf{x}) = \mathbf{1} \otimes \mathbf{p}(\mathbf{x})$,

$$\begin{aligned}
& \text{where } \mathbf{p}(\mathbf{x}) = \argmin_{\mathbf{y}_i} \mathbf{f}_i(\mathbf{x}_i) + \mathbf{h}_i(\mathbf{y}_i) \\
& \mathbf{x} \\
& \mathbf{y}_i \\
& \|\mathbf{x}_i\|_2 = \argmin_{\mathbf{y}_i} \|\mathbf{y}_i\|_2 \\
& \mathbf{1} \otimes \mathbf{x}_i \otimes \mathbf{1} \in \mathbb{R}^n \\
& (14)
\end{aligned}$$

Let $\mathbf{B} = \{\mathbf{x} : \|\mathbf{x}_i\|_2 \leq B, i = 1, \dots, N\}$. Hence, we can write $\mathbf{P}_{\mathbf{B}_0}(\mathbf{x}) = \mathbf{P}_{\mathbf{B}}(\mathbf{W} \otimes \mathbf{I}_n \mathbf{x})$ where $\mathbf{W} \in \mathbb{R}^{n \times n}$. Equivalently, $\mathbf{P}_{\mathbf{B}_0}(\mathbf{x}) = \mathbf{P}_{\mathbf{B}}(\mathbf{1} \otimes \mathbf{p}(\mathbf{x}))$, where $\mathbf{p}(\mathbf{x}) \in \mathbb{R}^n$. (15)

Although x -step and y -step of the PDA implementation in (12) can be computed locally at each node, computing \mathbf{y}_{k+1} requires communication among the nodes. Indeed, evaluating the average operator $\mathbf{p}(\cdot)$ is not a simple operation in a decentralized computational setting which only allows for communication

among neighbors. In order to overcome this issue, we will approximate $p(\cdot)$ operator using multi-consensus steps, and analyze the resulting iterations as an inexact primal-dual algorithm. In [20], this idea has been exploited within a distributed primal algorithm for unconstrained consensus optimization problems. We define the consensus step as one time exchanging local variables among neighboring nodes – the details of this operation will be discussed shortly. Since the connectivity network is dynamic, let $G_t = (N, E_t)$ be the connectivity network at the time t -th consensus step is realized for $t \in \mathbb{Z}^+$. We adopt the information exchange model in [21]. Assumption 3.1. Let $V_t \in \mathbb{R}^{N \times N}$ be the weight matrix corresponding to $G_t = (N, E_t)$ at the time of t -th consensus step and $N_{it} = \{j \in N : (i, j) \in E_t \text{ or } (j, i) \in E_t\}$. Suppose for all $t \in \mathbb{Z}^+$: (i) V_t is doubly stochastic; (ii) there exists $\alpha \in (0, 1)$ such that for $i \in N$, $V_{ijt} \geq \alpha$ if $j \in N_{it}$, and $V_{ijt} = 0$ if $j \notin N_{it}$; (iii) $G_t = (N, E_t)$ is connected where $E_t = \{(i, j) \in N \times N : t(i, j) \in E \text{ for infinitely many } t \in \mathbb{Z}^+\}$, and there exists $Z \in \mathbb{Z}^+$ such that if $(i, j) \in E$, then $(i, j) \in E_t$ for $t+1 \leq t \leq t+Z$ for all $t \geq 1$. Lemma 3.1. [21] Let Assumption 3.1 holds, and $W_{t,s} = V_t V_{t+1} \dots V_{s+1}$ for $t \leq s+1$. Given $s \geq 0$ the entries of $W_{t,s}$ converges to N^{-1} as $t \rightarrow \infty$ with a geometric rate, i.e., for all $i, j \in N$, one

$|W_{ijt,s} - N^{-1}| \leq \frac{2(1+\alpha)^T}{(1-\alpha)^T} \alpha^T$, where $\alpha = \frac{2(1+\alpha)^T}{(1-\alpha)^T}$, $\alpha \in (0, 1)$, $T = \frac{1}{\alpha}$, and $T \geq (N-1)T$.

Consider the k -th iteration of PDA as shown in (12). Instead of computing \bar{x}_{k+1} exactly according to (13), we propose to approximate \bar{x}_{k+1} with the help of Lemma 3.1 and set \bar{x}_{k+1} to this approximation. In particular, let t_k be the total number of consensus steps done before k -th iteration of PDA, and let $q_k \geq 1$ be the number of consensus steps within iteration k . For $x = [x_i]_{i \in N}$, define $R_k(x)$, $PB(W_{t_k+q_k, t_k} \bar{x})$

(16)

to approximate $PC(x)$ in (13). Note that $R_k(x)$ can be computed in a distributed fashion requiring q_k communications with the neighbors for each node. Indeed, $R_k(x) = [R_{ki}(x)]_{i \in N}$

such that $R_{ki}(x) = \frac{1}{N} \sum_{j \in N} W_{ijk, t_k+q_k} x_j$.

$\frac{1}{N} \sum_{j \in N}$

W_{ijk, t_k+q_k}

W_{ijk, t_k+q_k}

x_j .

(17)

Moreover, the approximation error, $R_k(x) - PC(x)$, for any x can be bounded as in (18) due to non-expansivity of PB and using Lemma 3.1. From (15), we get for all $i \in N$, $\|x_{t_k+q_k} - \bar{x}_{k+1}\|$

$\|R_{ki}(x) - \frac{1}{N} \sum_{j \in N} W_{ijk, t_k+q_k} x_j\| \leq \frac{1}{N} \sum_{j \in N} \|W_{ijk, t_k+q_k} x_j - \frac{1}{N} \sum_{j \in N} W_{ijk, t_k+q_k} x_j\|$

$\leq \frac{1}{N} \sum_{j \in N} \|W_{ijk, t_k+q_k} x_j - \frac{1}{N} \sum_{j \in N} W_{ijk, t_k+q_k} x_j\|$

$\leq \frac{1}{N} \sum_{j \in N} \|W_{ijk, t_k+q_k} x_j - \frac{1}{N} \sum_{j \in N} W_{ijk, t_k+q_k} x_j\|$

$\leq \frac{1}{N} \sum_{j \in N} \|W_{ijk, t_k+q_k} x_j - \frac{1}{N} \sum_{j \in N} W_{ijk, t_k+q_k} x_j\|$

$\leq \frac{1}{N} \sum_{j \in N} \|W_{ijk, t_k+q_k} x_j - \frac{1}{N} \sum_{j \in N} W_{ijk, t_k+q_k} x_j\|$

$\leq \frac{1}{N} \sum_{j \in N} \|W_{ijk, t_k+q_k} x_j - \frac{1}{N} \sum_{j \in N} W_{ijk, t_k+q_k} x_j\|$

$\leq \frac{1}{N} \sum_{j \in N} \|W_{ijk, t_k+q_k} x_j - \frac{1}{N} \sum_{j \in N} W_{ijk, t_k+q_k} x_j\|$

very slowly. We assume agents know q_k as a function of k at the beginning, hence, synchronicity can be achieved by simply counting local communications with each neighbor.

4

Numerical Section

We tested DPDA-S and DPDA-D on a primal linear SVM problem where the data is distributed among the computing nodes in N . For the static case, communication network $G = (N, E)$ is a connected graph that is generated by randomly adding edges to a spanning tree, generated uniformly at random, until a desired algebraic connectivity is achieved. For the dynamic case, for each consensus round $t \geq 1$, G_t is generated as in the static case, and $V_t, I_t \in \mathbb{R}^n$, where L_t is the Laplacian of G_t , and the constant $c \in \mathbb{R}$. We ran DPDA-S and DPDA-D on the line and complete graphs as well to see the topology effect for the dynamic case when the topology is line, each G_t is a random line graph. Let $S = \{1, 2, \dots, s\}$ and $D = \{(x_i, y_i) \in \mathbb{R}^n \times \{-1, +1\} : i \in S\}$ be a set of feature vector and label pairs. Suppose S is partitioned into S_{test} and S_{train} , i.e., the index sets for the test and training data; let $\{S_i\}_{i \in N}$ be a partition of S_{train} among the nodes N . Let $w = [w_i]_{i \in N}$, $b = [b_i]_{i \in N}$, and $w_i \in \mathbb{R} - S_{\text{train}}$ such that $w_i \in \mathbb{R}$ and $b_i \in \mathbb{R}$ for $i \in N$. Consider the following distributed SVM problem: $\min_{w, b, \gamma} \sum_{i \in N} \sum_{(x, y) \in S_i} \gamma (w^T x + b_i - y) \quad \text{s.t.} \quad \gamma \geq 0, \gamma = 0 \text{ if } S_i = \emptyset, w_i = w_j, b_i = b_j \text{ if } (i, j) \in E$

w, b, γ

$\sum_{i \in N} \sum_{(x, y) \in S_i} \gamma (w^T x + b_i - y)$

$\gamma \geq 0$

$w_i = w_j, b_i = b_j \text{ if } (i, j) \in E$

$\gamma = 0 \text{ if } S_i = \emptyset$

$\gamma (w^T x + b_i - y) \geq 0$

$\gamma \geq 0$

$\gamma (w^T x + b_i - y) \geq 0, \gamma = 0 \text{ if } S_i = \emptyset, w_i = w_j, b_i = b_j \text{ if } (i, j) \in E$

$\gamma \geq 0$

o

Similar to [3], $\{x_i\}_{i \in S}$ is generated from two-dimensional multivariate Gaussian distribution with covariance matrix $\Sigma = [1, 0; 0, 2]$ and with mean vector either $m_1 = [1, 1]^T$ or $m_2 = [1, -1]^T$ with equal probability. The experiment was performed for $C = 2, \gamma = 10, s = 900$ such that $|S_{\text{test}}| = 600, |S_i| = 30$ for $i \in N$, i.e., $|S_{\text{train}}| = 300$, and $q_k = k$. We ran DPDA-S and DPDA-D on line, random, and complete graphs, where the random graph is generated such that the algebraic connectivity is approximately 4. Relative suboptimality and relative consensus

violation, i.e., $\max_{(i, j) \in E} \|w_i - w_j\|_2 / \|w^* - b^*\|_2$, and absolute feasibility violation are

plotted against iteration counter in Fig. 3, where $\|w^* - b^*\|_2$ denotes the optimal solution to the central problem. As expected, the convergence is slower when the connectivity of the graph is weaker. Furthermore, visual comparison between DPDA-S, local SVMs (for two nodes) and centralized SVM for the same training and test data sets is given in Fig. 4 and Fig. 5 in the appendix.

Figure 3: Static (top) and Dynamic (bottom) network topologies: line, random, and complete graphs

8

2 References

- [1] Qing Ling and Zhi Tian. Decentralized sparse signal recovery for compressive sleeping wireless sensor networks. *Signal Processing, IEEE Transactions on*, 58(7):3816?3827, 2010.
- [2] Ioannis D Schizas, Alejandro Ribeiro, and Georgios B Giannakis. Consensus in ad hoc WSNs with noisy links - Part I: Distributed estimation of deterministic signals. *Signal Processing, IEEE Transactions on*, 56(1):350?364, 2008.
- [3] Pedro A Forero, Alfonso Cano, and Georgios B Giannakis. Consensus-based distributed support vector machines. *The Journal of Machine Learning Research*, 11:1663?1707, 2010.
- [4] Ryan McDonald, Keith Hall, and Gideon Mann. Distributed training strategies for the structured perceptron. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 456?464. Association for Computational Linguistics, 2010.
- [5] F. Yan, S. Sundaram, S. Vishwanathan, and Y. Qi. Distributed autonomous online learning: Regrets and intrinsic privacy-preserving properties. *Knowledge and Data Engineering, IEEE Transactions on*, 25(11):2483?2493, 2013.
- [6] Gonzalo Mateos, Juan Andr?es Bazerque, and Georgios B Giannakis. Distributed sparse linear regression. *Signal Processing, IEEE Transactions on*, 58(10):5262?5276, 2010.
- [7] Francis R Bach, Gert RG Lanckriet, and Michael I Jordan. Multiple kernel learning, conic duality, and the smo algorithm. In *Proceedings of the twenty-first international conference on Machine learning*, page 6. ACM, 2004.
- [8] Angelia Nedi?c and Asuman Ozdaglar. Subgradient methods for saddle-point problems. *Journal of optimization theory and applications*, 142(1):205?228, 2009.
- [9] Antonin Chambolle and Thomas Pock. A first-order primal-dual algorithm for convex problems with applications to imaging. *Journal of Mathematical Imaging and Vision*, 40(1):120?145, 2011.
- [10] Bingsheng He and Xiaoming Yuan. Convergence analysis of primal-dual algorithms for a saddle-point problem: from contraction perspective. *SIAM Journal on Imaging Sciences*, 5(1):119?149, 2012.
- [11] Antonin Chambolle and Thomas Pock. On the ergodic convergence rates of a first-order primal?dual algorithm. *Mathematical Programming*, 159(1):253?287, 2016.
- [12] Yunmei Chen, Guanghui Lan, and Yuyuan Ouyang. Optimal primal-dual methods for a class of saddle point problems. *SIAM Journal on Optimization*, 24(4):1779?1814, 2014.
- [13] A. Nedic and A. Ozdaglar. *Convex Optimization in Signal Processing and Communications*, chapter Cooperative Distributed Multi-agent Optimization, pages 340?385. Cambridge University Press, 2010.
- [14] A. Nedi?c. Distributed optimization. In *Encyclopedia of Systems and Control*, pages 1?12. Springer, 2014.
- [15] Tsung-Hui Chang, Angelia Nedic, and Anna Scaglione. Distributed constrained optimization by consensus-based primal-dual perturbation method. *Automatic Control, IEEE Transactions on*, 59(6):1524?1538, 2014.
- [16] David

Mateos-Núñez and Jorge Cortés. Distributed subgradient methods for saddle-point problems. In 2015 54th IEEE Conference on Decision and Control (CDC), pages 5462–5467, Dec 2015. [17] Deming Yuan, Shengyuan Xu, and Huanyu Zhao. Distributed primal-dual subgradient method for multi-agent optimization via consensus algorithms. *Systems, Man, and Cybernetics, Part B: Cybernetics*, IEEE Transactions on, 41(6):1715–1724, 2011. [18] Angelia Nedić, Asuman Ozdaglar, and Pablo A Parrilo. Constrained consensus and optimization in multiagent networks. *Automatic Control, IEEE Transactions on*, 55(4):922–938, 2010. [19] Kunal Srivastava, Angelia Nedić, and Dušan M Stipanović. Distributed constrained optimization over noisy networks. In *Decision and Control (CDC), 2010 49th IEEE Conference on*, pages 1945–1950. IEEE, 2010. [20] Albert I Chen and Asuman Ozdaglar. A fast distributed proximal-gradient method. In *Communication, Control, and Computing (Allerton), 2012 50th Annual Allerton Conference on*, pages 601–608. IEEE, 2012. [21] Angelia Nedić and Asuman Ozdaglar. Distributed subgradient methods for multi-agent optimization. *Automatic Control, IEEE Transactions on*, 54(1):48–61, 2009. [22] Ralph Tyrell Rockafellar. *Convex analysis*. Princeton university press, 2015. [23] H. Robbins and D. Siegmund. *Optimizing methods in statistics (Proc. Sympos., Ohio State Univ., Columbus, Ohio, 1971)*, chapter A convergence theorem for non negative almost supermartingales and some applications, pages 233 – 257. New York: Academic Press, 1971.