# Online Distributed Convex Optimization
# on Dynamic Networks

Saghar Hosseini, Airlie Chapman, and Mehran Mesbahi

*Abstract*—This paper presents a distributed optimization scheme over a network of agents in the presence of cost uncertainties and over switching communication topologies. Inspired by recent advances in distributed convex optimization, we propose a distributed algorithm based on dual sub-gradient averaging. A convergence rate analysis for the offline optimization, and a regret analysis for the online case, as a function of the underlying dynamic network topology are then presented for both classes of uncertainties. Application of the proposed setup is then discussed for uncertain sensor networks.

*Index Terms*—Distributed optimization; online optimization; switching graphs; weighted dual-averaging

## I. INTRODUCTION

In recent years, there has grown an extensive literature on distributed convex optimization on networks [2], [3], [4], and in particular, distributed subgradient algorithms [5], [6], [7]. The distributed nature of these algorithms is often dictated by constraints on the information flow amongst the nodes (agents) in the network, or mandated by the scale of the problem at hand. As such, distributed optimization aims to address the quality of the solution strategy necessitated by a *spatial* constraints, or uncertainties, due to the information flow across the network. Another class of uncertainties often encountered in applications is due to temporal variations in the cost structure and constraints. As such, robustness of the solution strategy with respect to temporal uncertainties becomes of paramount importance. One approach to improve the robustness of algorithms for convex optimization is via stochastic methods [8], [9], where the probability distribution of uncertain variable is known *a priori*. This approach has been pursued by Duchi *et al.* [6] that adopted a stochastic subgradient method where the distribution of subgradients is known *a priori*.

Despite its many successes, stochastic optimization-based methods do not explicitly address the dynamic aspect of the problem in an uncertain environment. Online learning is an extension of stochastic optimization where the uncertainty in the system is codified by an *arbitrarily* varying cost function revealed over time. In particular, at the time the relevant decision is made, the cost function is assumed to be unknown

and only revealed to the decision-maker after it commits to a decision. Such learning algorithms have had a significant impact on modern machine learning [10], [11]. One standard metric to measure the performance of these online algorithms is *regret*. Regret measures the difference between the incurred cost by the algorithm and the cost of the *best fixed decision in hindsight*. An online algorithm is then declared "good" when its regret is sub-linear. Distributed online optimization and its applications in multi-agent systems has not been studied at large by the systems and control community. Yan *et al.* in [12] introduced a decentralized online optimization based on the sub-gradient method in which the agents interact over a weighted strongly connected directed graph. Considering an undirected path graph with a fixed-radius neighborhood information structure, Raginsky *et al.* [13] proposed an online algorithm for distributed optimization based on sequential updates, proving a regret bound of $O(\sqrt{T})$. In [1], we proposed an extension to the work of Duchi *et al.* [6] on distributed optimization with convergence rate of $O(\sqrt{T} \log T)$ to an online setting. In addition, an improved regret bound of $O(\sqrt{T})$ has been derived for strongly connected networks, also highlighting the dependence of the regret on the connectivity of the underlying network. In this paper, we consider two classes of uncertainties arising in distributed convex optimization, namely, uncertainties in the cost and the network structure. Our main contribution is to show that distributed dual averaging provides an effective mean of dealing with these classes of uncertainties for large-scale decision-making.

The organization of the paper is as follows. The notation and background on graphs and regret are reviewed in §II. In §III, the formulation of distributed convex optimization and its solution via dual averaging are presented. The convergence analysis for distributed convex optimization via dual averaging is then extended to switching topologies in §IV-A. In §IV-B, distributed convex optimization problem is extended to the online setting. In §IV-C, we discuss the application of the proposed framework for online distributed estimation over uncertain sensor networks.

## II. BACKGROUND AND PRELIMINARIES

We provide a brief background on constructs that will be used in this paper. For the column vector $v \in \mathbb{R}^p$, $v_i$ or $[v]_i$ denotes its $i$th element. For matrix $M \in \mathbb{R}^{p \times q}$, $[M]_{ij}$, or simply $M_{ij}$, denotes the element in its $i$th row and $j$th column. The family of probability vectors is denoted by $\Omega$ and contains all non-negative vectors $\sigma \in [0,1]^n$ such that

$\sum \sigma_i = 1$. A row stochastic matrix $P$ is a non-negative matrix with rows in $\Omega$; the ergodic coefficient for a stochastic matrix $Q \in \mathbb{R}^{n \times n}$ is denoted by

$$\tau(Q) = 1 - \min_{i,j \in [n]} \sum_{k=1}^{n} \min\{Q_{ik}, Q_{jk}\}. \qquad (1)$$

A time varying matrix is denoted by $P^t$; a (backward) sequence of time varying matrices on the other hand is designated by $P^{(t,0)} = P^t P^{t-1} \cdots P^0$. For any positive integer $n$, the set $\{1, 2, ..., n\}$ is denoted by $[n]$. The inner product of two vectors $\theta$ and $\phi$ is represented by $\langle \theta, \phi \rangle$. The 2-norm is signified by $||.||_2$, a general norm of vector is denoted as $||\theta||$, and its associated dual norm is defined as $||\theta||_* = \sup_{||\phi||=1} \langle \theta, \phi \rangle$. A function $f : \Theta \to \mathbb{R}$, where $\Theta \subseteq \mathbb{R}^m$ for some positive integer $m$, is called $L$-Lipschitz continuous with respect to the norm $||\cdot||$ if there exists a positive constant $L$ for which

$$|f(\theta) - f(\phi)| \le L \|\theta - \phi\| \text{ for all } \theta, \phi \in \Theta. \qquad (2)$$

Although the dual of the 2-norm is the 2-norm itself, we derive some of the bounds in our subsequent analysis for a more general setting using the notion of the dual norm. The main reason is the connection between the Lipschitz continuity of a function (in the native norm) and the boundedness of its subgradient (by the Lipschitz constant) in the dual norm.

Regret is one measure of the performance for learning algorithms. In the online optimization setting, an algorithm generates a sequence of decisions $\{x(t)\}_{t=1}^{T}$. The number of iterations $T$ is unknown to the online algorithm. At each iteration $t$, after committing to $x(t)$, a previously unknown convex cost function $f_t$ is revealed, and a loss $f_t(x(t))$ is incurred by the algorithm. The goal of the online algorithm is to ensure that the time average of the difference between the cost incurred by the algorithm and the cost of the best fixed decision $x^* = \text{argmin} \sum_{t=1}^{T} f_t(x)$ is small. This difference, namely,

$$R_T(x^*, x) = \sum_{t=1}^{T} (f_t(x(t)) - f_t(x^*)), \qquad (3)$$

is referred to as the regret of the online algorithm. An algorithm performs well if its regret is sub-linear as a function of $T$, i.e. $\lim_{T \to \infty} R_T/T = 0$. This implies that on average, the algorithm performs as well as the best fixed strategy in hindsight independent of the adversary's moves and environmental uncertainties [14], [15]. In order to analyze the performance of *distributed online algorithms* two variations of the notion of regret will be considered. First, is the regret due to agent $i$'s decision,

$$R_T(x^*, x_i) = \sum_{t=1}^{T} (f_t(x_i(t)) - f_t(x^*)), \qquad (4)$$

and the second, is the regret based on the (temporal) *running average* of the decisions, namely, $\widetilde{x}_i = \frac{1}{T} \sum_{t=1}^{T} x_i(t)$,

$$R_T(x^*, \widetilde{x}_i) = \sum_{t=1}^{T} (f_t(\widetilde{x}_i(t)) - f_t(x^*)). \qquad (5)$$

## III. PROBLEM SETUP

In this section a distributed decision process is considered in which a network of agents cooperatively optimize a global objective function. The objective to be minimized is

$$f(x) = \frac{1}{n} \sum_{i=1}^{n} f_i(x) \quad \text{subject to } x \in \chi; \qquad (6)$$

the following standing assumptions are in place: (**a**) the local cost function $f_i : \mathbb{R}^d \to \mathbb{R}$ associated with agent $i \in V$ is convex for all $i$, (**b**) the decision space $\chi \subseteq \mathbb{R}^d$ is a closed convex set, (**c**) each convex function $f_i$ is $L$-Lipschitz with respect to $||.||$, (**d**) the underlying network $\mathcal{G}$ is strongly connected.[1] Furthermore, for switching topologies we assume that, (**e**) the union of directed topologies $\mathcal{G}^{\cup_{i=0}^{\delta-1}} = \bigcup_{i=0}^{\delta-1} \mathcal{G}^i$ over some fixed uniform interval of length $\delta$, with $\delta \ge 1$ a positive integer, is strongly connected.

In order to solve the optimization problem (6) over switching networks and in an online setting, we adapt Nesterov's dual averaging algorithm [17] and our preliminary results [1], that in turn, has been inspired by [6].

The centralized form of the dual averaging algorithm appears as a sub-gradient descent method followed by a projection step onto the constraint set $\chi$, specifically we let,

$$y(t+1) = y(t) + g(t), \qquad (7)$$

where $g(t) \in \partial f(x(t))$ followed by the primal update,

$$x(t+1) = \Pi_{\chi}^{\psi}(y(t+1), \alpha(t)), \qquad (8)$$

where $\Pi_{\chi}^{\psi}(\cdot)$ is a regularized projection onto $\chi$. The projection is represented by

$$\Pi_{\chi}^{\psi}(y(t), \alpha(t)) = \arg \min_{x \in \chi} \left\{ \langle y(t), x \rangle + \frac{1}{\alpha(t)} \psi(x) \right\}, \qquad (9)$$

where $\alpha(t)$ is a non-increasing sequence of positive functions and $\psi(x) : \chi \to \mathbb{R}$ is a proximal function. The standard dual averaging algorithm uses the proximal function $\psi(x)$ to avoid undesirable oscillations in the projection step. Without loss of generality, $\psi$ is assumed to be strongly convex with respect to $||.||$, $\psi(x) \ge 0$, and $\psi(0) = 0$.

An extension of dual averaging to the network setting is presented in Algorithm 1. This so-called Distributed Weighted Dual Averaging (DWDA) algorithm sequentially updates the *local* $x_i(t)$ and a *working* variable $y_i(t)$ for each agent $i$ at time $t$. The update itself is based on a local subgradient of the cost $f_i(x_i(t))$ denoted by $g_i(t)$. This distributed algorithm can be considered as an approximate subgradient descent; the approximation is attained by an agent via a convex combination of local subgradients provided by its

---

[1]A directed graph is strongly connected if there is a directed path between every pair of nodes [16].

neighbors. This operation can be represented compactly by a stochastic matrix $P \in \mathbb{R}^{n \times n}$.

---

**Algorithm 1:** Distributed Weighted Dual Averaging (DWDA)

1  **for** $t = 1$ **to** $T$ **do**
2      Evaluate $f(t) = \{f_i(t); \text{ for all } i = 1, ..., n\}$
3      **foreach** *Agent i* **do**
4          Compute subgradient $g_i(t) \in \partial f_i(x_i(t))$
5          $y_i(t+1) = \sum_{j \in N(i)} P_{ji}(t) y_j(t) + g_i(t)$
6          $x_i(t+1) = \prod_{\chi}^{\psi}(y_i(t+1), \alpha(t))$
7          $\widetilde{x}_i(t+1) = \frac{1}{t+1}\sum_{s=1}^{t+1} x_i(s)$
8      **end**
9  **end**

---

Before presenting the convergence analysis of the distributed optimization algorithm for the online setting and switching networks, it is instructive to provide a few preliminary remarks on the convergence analysis in the offline setting for a fixed strongly connected network. First, since this fixed network is strongly connected, the communication matrix $P$ is 1-irreducible ([18]; Corollary 4); in fact, as the diagonal elements are positive, this matrix is indecomposable and aperiodic (SIA). Now, given that $P$ is SIA, there exists a vector $\pi \in \Omega$ [19], such that

$$\pi_j = \sum_{i=1}^{n} \pi_i P_{ij}. \tag{10}$$

We thereby consider the sequences $\bar{y}(t)$ and $\bar{g}(t)$ defined as

$$\bar{y}(t) = \sum_{i=1}^{n} \pi_i y_i(t), \text{ and } \bar{g}(t) = \sum_{i=1}^{n} \pi_i g_i(t), \tag{11}$$

signifying the (network-level) weighted average of dual variables and subgradients in the DWDA algorithm, respectively. Based on (10) and (11) it follows that

$$
\begin{aligned}
\bar{y}(t+1) &= \sum_{i=1}^{n} \pi_i \left\{ \sum_{j=1}^{n} P_{ji} y_j(t) + g_i(t) \right\} \\
&= \sum_{j=1}^{n} y_j(t) \pi_j + \bar{g}(t) = \bar{y}(t) + \bar{g}(t),
\end{aligned} \tag{12}
$$

which is analogous to the dual averaging update (7). Thus, the following update rule is introduced which is analogous to the standard dual averaging algorithm projection step (8), now on the network-level weighted average of dual variables,

$$\phi(t+1) = \Pi_{\chi}^{\psi}(\bar{y}(t+1), \alpha(t)). \tag{13}$$

Using these observations, Duchi *et al.* [6] showed that given the sequences $x_i(t)$ and $y_i(t)$ generated by lines 5 and 6 of Algorithm 1, for all $i \in [n]$,

$$
\begin{aligned}
\frac{1}{T}\sum_{t=1}^{T} f(x_i(t)) - f(x^*) \leq\ & \frac{L^2}{2}\sum_{t=1}^{T}\alpha(t-1) + \frac{1}{\alpha(T)}\psi(x^*) \\
& + \frac{L}{T}\sum_{t=1}^{T}\alpha(t)(\|\bar{y}(t) - y_i(t)\|_*) \\
& + \frac{2}{n}\sum_{i=1}^{n}\|\bar{y}(t) - y_i(t)\|_*). 
\end{aligned} \tag{14}
$$

The last two terms on the right hand side of (14) represent the error due to the network which is defined as the deviation of local dual variable $y_i$ from the weighted average of the dual variable $\bar{y}$ over the network. As such, a bound on this term for fixed or switching network proves to be essential for convergence rate as well as regret analysis.

## IV. MAIN RESULTS

We now consider the extensions of the DWDA algorithm for the situations where the network is allowed to switch and the cost function has an online character. We then proceed to apply the setup to a distributed estimation scenario on uncertain sensor networks.

### A. Switching Networks

The main idea behind embedding switching networks in distributed decision-making and reasoning about convergence properties, as previously used in the literature, e.g., [20], is the weak ergodicity of inhomogeneous Markov chains. In this direction, enforcing assumption (**e**) becomes particularly pertinent. This assumption implies that the product of stochastic SIA matrices associated with the underlying switching network converges exponentially to a rank-one matrix of the form $\mathbf{1}\pi^T$ as $t \to \infty$, where $\pi \in \Omega$. In fact, analogous to the analysis used in [20], it follows that the product

$$P^{(k\delta-1,(k-1)\delta)} \cdots P^{(2\delta-1,\delta)} P^{(\delta-1,0)}$$

converges exponentially to a rank-one matrix of the form $\mathbf{1}\pi^T$ as $t \to \infty$; moreover, based on Theorem 1 of [19], we have

$$\left| P^{(k\delta-1,0)} - \pi_j \right| \leq \gamma^{\left\lfloor \frac{k}{\nu} \right\rfloor}, \tag{15}$$

where

$$\gamma = \max_{\nu \geq 1}\left\{ \tau(P^{(\delta\nu-1,0)}) < 1 \right\}; \tag{16}$$

note that the maximization is over all realizations of the sequence $P^{(\delta\nu-1,0)}$; the parameter $\nu$ can be shown to be bounded by the diameter of the network.[2]

Before we state the result on the rate of convergence of DWDA over switching graphs, a key observation is in order. This result imposes an upper bound on the effect of the network topology associated with $\|\bar{y}(t) - y_i(t)\|_*$ in (14) proportional to the error incurred by the decentralized update in Algorithm 1; moreover this bound highlights the importance of the underlying network topology through the communication matrix $P^t$ and its products.

**Lemma 1.** *For all $i \in [n]$ the sequences $y_i(t)$ and $\bar{y}(t)$ generated by line 5 of Algorithm 1 and (12), satisfy,*

$$\|\bar{y}(t) - y_i(t)\|_* \leq L \sum_{k=0}^{t-2}\sum_{j=1}^{n}\left| P_{ij}^{(t-1,k+1)} - \pi_j \right| + 2L.$$

---

[2] The maximum path length between any pair of nodes in the network.

*Proof:* Reformulating this update and by induction through $s$ steps we have,

$$y_i(t) = \sum_{j=1}^{n} P_{ij}^{(t-1,t-s)} y_j(t-s)$$
$$+ \sum_{k=t-s}^{t-2} \sum_{j=1}^{n} P_{ij}^{(t-1,k+1)} g_j(k) + g_i(t-1). \quad (17)$$

Since $\bar{y}(t)$ evolves as in (12), by setting $s = t$ in (17) and assuming $y_i(0) = 0$, we get,

$$\bar{y}(t) - y_i(t) = \sum_{k=0}^{t-2} \left( \sum_{j=1}^{n} \left( \pi_j - P_{ij}^{(t-1,k+1)} \right) g_j(k) \right)$$
$$+ \bar{g}(t-1) - g_i(t-1). \quad (18)$$

Thus, the dual norm of (18) is bounded as

$$\|\bar{y}(t) - y_i(t)\|_* \le \| \sum_{k=0}^{t-2} \left( \sum_{j=1}^{n} \left( \pi_j - P_{ij}^{(t-1,k+1)} \right) g_j(k) \right) \|_*$$
$$+ \|\bar{g}(t-1) - g_i(t-1)\|_*;$$
$$\le \sum_{k=0}^{t-2} \sum_{j=1}^{n} \left| P_{ij}^{(t-1,k+1)} - \pi_j \right| \|g_j(k)\|_*$$
$$+ \|\bar{g}(t-1) - g_i(t-1)\|_*.$$

Since $\|g_i(t)\|_\star \le L$,

$$\|\bar{y}(t) - y_i(t)\|_* \le L \sum_{k=0}^{t-2} \sum_{j=1}^{n} \left| P_{ij}^{(t-1,k+1)} - \pi_j \right| + 2L. \quad (19)$$

∎

We are ready to state the convergence result for distributed optimization via dual averaging on switching networks.

**Theorem 2.** *Given the sequences $x_i(t)$ and $y_i(t)$ generated by lines 5 and 6 in* Algorithm 1, *for all $i \in [n]$ with $\psi(x^*) \le R^2$ and $\alpha(t) = k/\sqrt{t}$, we have*

$$\frac{1}{T} \sum_{t=1}^{T} f(x_i(t)) - f(x^*)$$
$$\le \left( \frac{R^2}{k} + kL^2 \left( \frac{6n}{1-\gamma} + 6n\delta\nu + 1 \right) \right) \frac{1}{\sqrt{T}}, \quad (20)$$

*where $\gamma < 1$ is a function of the ergodicity of the communication matrix (see (16)) while $\nu$ is a measure of network connectivity and is bounded by the diameter of the network.*

*Proof:* Based on (15) and Lemma 1, it follows that

$$\|\bar{y}(t) - y_i(t)\|_* \le nL \sum_{k=1}^{t-1} \gamma^k + nL(\delta\nu - 1) + 2L, \quad (21)$$

and since $\gamma < 1$, (21) is further bounded by,

$$\|\bar{y}(t) - y_i(t)\|_* \le nL \left( \frac{1}{1-\gamma} + \delta\nu - 1 \right) + 2L. \quad (22)$$

Therefore, the integral test on $\alpha(t) = k/\sqrt{t}$ provides the

bound on the first and last terms in (20) as[3]

$$\frac{1}{T} \sum_{t=1}^{T} f(x_i(t)) - f(x^*) \le \frac{kL^2}{\sqrt{T}} + \frac{\psi(\theta^\star)}{k\sqrt{T}}$$
$$+ \frac{6kL^2}{\sqrt{T}} \left( \frac{n}{1-\gamma} + n\delta\nu + 1 \right).$$

Given $\psi(x^*) \le R^2$, the statement of the theorem now follows. ∎

Theorem 2 states that Algorithm 1 performs "well" as it exhibits a sub-linear convergence rate. It also highlights the importance of the underlying network topology through the parameters $\gamma$ and $\nu$.

### B. Online Distributed Optimization on Switching Networks

We now consider the effect of uncertainties in the environment on distributed decision-making process where the global objective is to minimize

$$f_t(x) = \frac{1}{n} \sum_{i=1}^{n} f_{t,i}(x) \quad \text{subject to } x \in \chi, \quad (23)$$

where $f_{t,i} : \mathbb{R}^d \to \mathbb{R}$ is a convex cost function associated with agent $i \in V$, assumed to be revealed to the agent only after the agent commits to its decision. In other words, the function $f_{t,i}$ is allowed to change over time in an unpredictable manner due to modeling errors and uncertainties in the environment. The optimization variable $x_i \in \mathbb{R}^d$ belongs to a closed convex set $\chi \subseteq \mathbb{R}^d$ and represents the local decision made by agent $i$. The regret analysis presented below is for the *online* version of the DWDA presented in Algorithm 1.

**Theorem 3.** *Given the sequences $x_i(t)$ and $y_i(t)$ generated by lines 5 and 6 in* Algorithm 1, *for all $i \in [n]$ with $\psi(x^*) \le R^2$ and $\alpha(t) = k/\sqrt{t}$, we have*

$$R_T(x^*, x_i) \le \left( \frac{R^2}{k} + kL^2 \left( \frac{6n}{1-\gamma} + 6n\delta\nu + 1 \right) \right) \sqrt{T}, \quad (24)$$

*where $\gamma$ is a function of the ergodicity of the communication matrix (see (16)) while $\nu$ is a measure of network connectivity and is bounded by the diameter of the network.*

*Proof:* Consider an arbitrary fixed decision $x^* \in \chi$ and a sequence $\phi(t)$ generated by (13). From the $L$-Lipschitz continuity of $f_{t,i}$'s and the definition of regret (4), we have

$$R_T(x^*, x_i) \le \sum_{t=1}^{T} \left( f_t(\phi(t)) - f_t(x^*) + L\|x_i(t) - \phi(t)\| \right). \quad (25)$$

Note that we can reformulate the first term on the right hand side of (25) as

$$f_t(\phi(t)) - f_t(x^*) = \left( \frac{1}{n} \sum_{i=1}^{n} f_{t,i}(x_i(t)) - f_t(x^*) \right)$$
$$+ \left( \frac{1}{n} \sum_{i=1}^{n} [f_{t,i}(\phi(t)) - f_{t,i}(x_i(t))] \right). \quad (26)$$

---

[3]Note that $\sum_{t=1}^{T} \frac{k}{\sqrt{t}} \le 2k\sqrt{T} - k$.

Based on the convexity of $f_{t,i}$'s, we have

$$\sum_{t=1}^{T} \left( \frac{1}{n} \sum_{i=1}^{n} f_{t,i}(x_i(t)) - f_t(x^*) \right)$$
$$\leq \sum_{t=1}^{T} \left( \frac{1}{n} \sum_{i=1}^{n} \langle g_i(t), x_i(t) - x^* \rangle \right), \qquad (27)$$

where $g_i(t) \in \partial f_{t,i}(x_i(t))$ is the sub-gradient of $f_{t,i}$ at $x_i(t)$. Thereby, we can express the regret bound based on (26), (27), and the $L$-Lipschitz continuity of $f_{t,i}$'s as,

$$R_T(x^*, x_i) \leq \sum_{t=1}^{T} \left( \frac{1}{n} \sum_{i=1}^{n} \langle g_i(t), x_i(t) - x^* \rangle \right.$$
$$\left. + \frac{L}{n} \sum_{i=1}^{n} \|x_i(t) - \phi(t)\| + L\|x_i(t) - \phi(t)\| \right). \qquad (28)$$

The first term on the right had side of (28) can be expanded as

$$\sum_{t=1}^{T} \left( \frac{1}{n} \sum_{i=1}^{n} \langle g_i(t), x_i(t) - x^* \rangle \right) \qquad (29)$$

$$= \sum_{t=1}^{T} \left( \frac{1}{n} \sum_{i=1}^{n} \langle g_i(t), x_i(t) - \phi(t) \rangle + \frac{1}{n} \sum_{i=1}^{n} \langle g_i(t), \phi(t) - x^* \rangle \right). \qquad (30)$$

Now, we need to bound the terms on the right hand side of (30). The first term is bounded based on the convexity and $L$-Lipschitz continuity of $f_{t,i}$.[4] In other words,

$$\langle g_i(t), x_i(t) - \phi(t) \rangle \leq L\|x_i(t) - \phi(t)\|. \qquad (31)$$

Since $x_i(t)$ and $\phi(t)$ are the projections of $y_i(t)$ and $\bar{y}(t)$ respectively, the Lipschitz continuity of $\Pi_{\chi}^{\psi}(., \alpha)$ (see the Appendix) imposes a bound on $\|x_i(t) - \phi(t)\|$ as

$$\|x_i(t) - \phi(t)\| \leq \alpha(t)\|\bar{y}(t) - y_i(t)\|_*, \qquad (32)$$

where $\|.\|_*$ is the dual norm. Noting that $\|g_i(t)\|_* \leq L$, we can thus write

$$\sum_{t=1}^{T} \left( \frac{1}{n} \sum_{i=1}^{n} \langle g_i(t), x_i(t) - x^* \rangle \right)$$
$$\leq \frac{L}{n} \sum_{t=1}^{T} \sum_{i=1}^{n} \alpha(t)\|\bar{y}(t) - y_i(t)\|_*$$
$$+ \frac{L^2}{2} \sum_{t=2}^{T} \alpha(t-1) + \frac{1}{\alpha(T)} \psi(x^*). \qquad (33)$$

Thereby, (28), (32), and (33) imply that

$$R_T(x^*, x_i) \leq \frac{L^2}{2} \sum_{t=2}^{T} \alpha(t-1) + \frac{1}{\alpha(T)} \psi(x^\star)$$
$$+ L \sum_{t=1}^{T} \alpha(t) \left( \|\bar{y}(t) - y_i(t)\|_* + \frac{2}{n} \sum_{i=1}^{n} \|\bar{y}(t) - y_i(t)\|_* \right)$$

On the other hand, Lemma 1 imposes an upper bound on the last term on the right hand side of the above in equality. Thus, using (22), the regret is further bounded as

$$R_T(x^*, x_i) \leq \frac{L^2}{2} \sum_{t=1}^{T-1} \alpha(t) + \frac{1}{\alpha(T)} \psi(x^*)$$
$$+ 3L^2 \left( \frac{n}{1-\gamma} + n\delta\nu + 2(1-n) \right) \sum_{t=1}^{T} \alpha(t). \qquad (34)$$

---

[4]Note that convexity of $f_{t,i}$ implies $\langle g_i(t), x - y \rangle \leq f_{t,i}(x) - f_{t,i}(y)$. Therefore, based on $L$-Lipschitz continuity of $f_{t,i}$'s, we have $\|g_i\|_* \leq L$ and we can deduce (31).

The statement of the theorem now follows from the integral test on $\alpha(t) = k/\sqrt{t}$ and $\psi(x^*) \leq R^2$. ∎

Theorem 3 indicates "good" performance of online-DWDA through sub-linear regret and highlights the importance of the underlying network topology through the parameters $\gamma$ and $\nu$ examined in §IV. The regret analysis for the (temporal) running average estimates for each agent exhibit a similar dependence on the network structure.

**Corollary 4.** *Given the sequence $\widetilde{x}_i(t)$ generated by line* 7 *in* Algorithm 1 *with $\psi(x^*) \leq R^2$ and $\alpha(t) = k/\sqrt{t}$, we have*

$$R_T(x^*, \widetilde{x}_i) \leq 2 \left( \frac{R^2}{k} + kL^2 \left( \frac{6n}{1-\gamma} + 6n\delta\nu + 1 \right) \right) \sqrt{T}.$$

*Proof:* Since the cost function $f_t(x(t))$ is convex, $f_t(\widetilde{x}_i(t)) \leq \frac{1}{t} \sum_{s=1}^{t} f_t(x_i(s))$. Therefore, we have

$$f_t(\widetilde{x}_i(t)) - f_t(x^*) \leq \frac{1}{t} R_t(x^*, x_i). \qquad (35)$$

Thus, the running average regret is bounded as

$$R_T(x^*, \widetilde{x}_i) \leq \sum_{t=1}^{T} \left( \frac{1}{t} R_t(x^*, x_i) \right). \qquad (36)$$

In the meantime, the regret bound (24) in conjunction with the integral test implies the statement of the Corollary. ∎

### C. Example: Online Distributed Estimation

We consider a distributed sensor network that collectively aim to estimate a random vector

$$\theta \in \Theta = \left\{ \theta \in \mathbb{R}^d | \|\theta\|_2 \leq \theta_{\max} \right\}$$

at time $t$. The observation vector $z_{t,i} : \mathbb{R}^d \to \mathbb{R}^{p_i}$ represents the $i$th sensor measurement at time $t$ which is uncertain and time-varying due to the sensor's susceptibility to unknown environmental factors such as jamming. The sensor is assumed (not necessarily accurately) to have a linear model of the form $h_i(\theta) = H_i\theta$, where $H_i \in \mathbb{R}^{p_i \times d}$ is the observation matrix of sensor $i$ and $\|H_i\|_1 \leq h_{\max}$ for all $i$. The objective is to find the argument $\hat{\theta}$ that minimizes the cost function

$$f_t(\hat{\theta}) = \frac{1}{n} \sum_{i=1}^{n} f_{t,i}(\hat{\theta}) \quad \text{subject to } \hat{\theta} \in \Theta, \qquad (37)$$

where $f_{t,i}(\hat{\theta}) = \frac{1}{2} \left\| z_{i,t} - H_i\hat{\theta} \right\|_2^2$ is a convex cost function associated with sensor $i \in [n]$. It is assumed that the value of this local cost at time $t$ is only revealed to the sensor after $\hat{\theta}(t)$ has been computed, that is, the local error functions are allowed to change over time in an uncertain manner due to modeling errors and uncertainties in the environment. An online framework is particularly suitable for estimation problems without relying on prior assumption or knowledge of the statistical properties of the data. In the proposed distributed estimation algorithm, at time step $t$, each sensor $i$ estimates $\hat{\theta}_i \in \Theta$ based on the local information available to it and then an "oracle" reveals the cost $f_t(\hat{\theta}_i)$. The bounds presented in Theorem 3 apply after selecting $\psi(\hat{\theta}) = \frac{1}{2}\|\hat{\theta}\|_2^2$ and the parameter $\alpha(t)$ accordingly. In order to find the
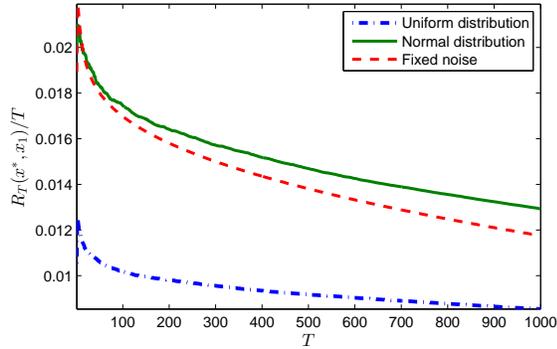
Fig. 1. Regret performance for three different observation noise characteristics, where $\mathcal{G}$ is a 100 node random 4-regular graph. The noise signals have been generated from distributions with mean $-b_{\max}$ and standard deviation $b_{\max}$.

constants $R$ and $L$ featured in the result, we note that for $\hat{\theta} \in \Theta$, $\psi(\hat{\theta}) \le (1/2)\theta_{\max}^2$, and thus $R \le 1/\sqrt{2}\theta_{\max}$. In this example, we assume that the observation for agent $i$ at time $t$ is of the form $z_{t,i} = a_t\theta + b_t$ for some $a \in (0, a_{\max})$ and $b \in (-b_{\max}, b_{\max})$. Therefore,

$$\sup_{\theta \in \chi} \|z_{t,i}(\theta)\|_2 \le a_{\max}\theta_{\max} + b_{\max}.$$

Further, the function $f_{t,i}$ is Lipschitz as it is convex on a compact domain; the Lipschitz constant can be shown to be $L = ((1/2)\theta_{\max}h_{\max} + a_{\max}\theta_{\max} + b_{\max})\, h_{\max}$. Hence $R_T(\theta^*, \hat{\theta}_i)/T \to 0$ and the algorithm performs as well as best fixed estimate $\theta^*$ in hindsight "on average". For the case where $\theta_t = \theta_{t+1}$ for $t = 1, 2, \ldots, T$, $\theta^*$ is the optimal estimate. The DWDA algorithms have been implemented on the described distributed sensor setup for $n = 100$ sensors over a random network with edge probability $p = 0.08$. The objective is to estimate a scalar $\theta \in (-1/2, 1/2)$ with a fixed $H_i \in (0, 1/4)$ for each agent; hence $\sup_i |H_i| = 1/4$. In this example, we have assumed $a \in (0, 1)$, $b \in (-1/4, 1/4)$, $\beta = 0.9$, and $k = 1/4$. Thus, $d = 1$, $\Theta = (-1/2, 1/2)$, $h_{\max} = 1/4$, $\theta_{\max} = 1/2$, $R = 1/(2\sqrt{2})$, and $L = 13/64$. The performance of the proposed online distributed estimation in the presence of various noise types is presented in Figure 1. These simulation results indicate that $R_T(\theta^*, \hat{\theta}_1) = O(\sqrt{T})$ for all noise types considered without a prior assumption on the noise characteristics.

## V. Conclusion

This paper studies the problem of decentralized optimization on dynamic networks operating in an uncertain environment. The uncertainties considers are due to network structure as well as the cost structure. Our approach indicates that dual averaging provides an effect means of dealing with both classes of uncertainties in a distributed setting, while also providing a privacy feature, where primal variable are not shared on the network during the optimization process.

## VI. Appendix

The following two observations shown by Duchi *et al.*, [6] have been used in the paper: (1) For any $u, v \in \mathbb{R}^m$, and under the conditions stated for proximal function $\psi$ and step size $\alpha(t)$, we have $\|\Pi_\chi^\psi(u, \alpha) - \Pi_\chi^\psi(v, \alpha)\| \le \alpha\|u - v\|_*$, and (2) For any positive and non-increasing sequence $\alpha(t)$ and $x^* \in \chi$,

$$\sum_{t=1}^{T} \langle \bar{g}(t), \phi(t) - x^*(t) \rangle \le \frac{1}{2}\sum_{t=1}^{T} \alpha(t-1)\|\bar{g}(t)\|_\star^2 + \frac{1}{\alpha(T)}\psi(x^*),$$

where the sequence $\phi(t)$ is generated by (13).

## References

[1] S. Hosseini, A. Chapman, and M. Mesbahi, "Online Distributed Optimization via Dual Averaging," in *IEEE Conference on Decision and Control*, 2013, pp. 1484 – 1489.

[2] S. Boyd, "Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers," *Foundations and Trends in Machine Learning*, vol. 3, no. 1, pp. 1–122, 2010.

[3] D. Jakovetic, J. Xavier, and J. M. F. Moura, "Cooperative Convex Optimization in Networked Systems: Augmented Lagrangian Algorithms With Directed Gossip Communication," *IEEE Transactions on Signal Processing*, vol. 59, no. 8, pp. 3889–3902, Aug. 2011.

[4] D. Mosk-Aoyama, T. Roughgarden, and D. Shah, "Fully distributed algorithms for convex optimization problems," *Distributed Computing*, vol. 4731, pp. 492–493, 2007.

[5] A. Nedic and A. Ozdaglar, "Distributed subgradient methods for multi-agent optimization," *IEEE Transactions on Automatic Control,*, vol. 54, pp. 48–61, 2009.

[6] J. C. Duchi, A. Agarwal, and M. J. Wainwright, "Dual averaging for distributed optimization: convergence analysis and network scaling," *IEEE Transactions on Automatic Control*, vol. 57, no. 3, pp. 592–606, 2012.

[7] S. Lee and A. Nedic, "Distributed random projection algorithm for convex optimization," *IEEE Journal of Selected Topics in Signal Processing*, vol. 7, pp. 221–229, 2013.

[8] S. Sundhar Ram, A. Nedić, and V. V. Veeravalli, "Distributed Stochastic Subgradient Projection Algorithms for Convex Optimization," *Journal of Optimization Theory and Applications*, vol. 147, no. 3, pp. 516–545, 2010.

[9] A. Agarwal and J. Duchi, "Distributed delayed stochastic optimization," in *IEEE Conference on Decision and Control*, 2012, pp. 5451–5452.

[10] M. Zinkevich, "Online convex programming and generalized infinitesimal gradient ascent," in *International Conference on Machine Learning*, 2003, pp. 421–422.

[11] L. Xiao, "Dual Averaging Methods for Regularized Stochastic Learning and Online Optimization," *Journal of Machine Learning Research*, vol. 11, pp. 2543–2596, 2010.

[12] F. Yan, S. Sundaram, S. V. N. Vishwanathan, and Y. Qi, "Distributed autonomous online learning: Regrets and intrinsic privacy-preserving properties," *IEEE Transactions on Knowledge and Data Engineering*, vol. 25, pp. 2483 – 2493, 2013.

[13] M. Raginsky, N. Kiarashi, and R. Willett, "Decentralized Online Convex Programming with Local Information," in *American Control Conference*, 2011, pp. 5363–5369.

[14] S. Shalev-Shwartz, "Online learning and online convex optimization," *Foundations and Trends in Machine Learning*, vol. 4, pp. 107–194, 2012.

[15] E. Hazan, "The Convex Optimization Approach to Regret Minimization," *Optimization for machine learning*, pp. 287–294, 2011.

[16] M. Mesbahi and M. Egerstedt, *Graph Theoretic Methods in Multiagent Networks*. NJ: Princeton University Press, 2010.

[17] Y. Nesterov, "Primal-dual subgradient methods for convex problems," *Mathematical Programming*, vol. 120, pp. 221–259, 2007.

[18] C. W. Wu, "On bounds of extremal eigenvalues of irreducible and m-reducible matrices," *Linear Algebra and its Applications*, vol. 402, pp. 29–45, 2005.

[19] J. Anthonisse and H. Tijms, "Exponential convergence of products of stochastic matrices," *Journal of Mathematical Analysis and Applications*, vol. 59, no. 2, pp. 360–364, 1977.

[20] A. Jadbabaie and A. Morse, "Coordination of groups of mobile autonomous agents using nearest neighbor rules," *IEEE Transactions on Automatic Control*, vol. 48, no. 6, pp. 988–1001, 2003.