# LOGARITHMIC REGRET BOUND OVER DIFFUSION BASED DISTRIBUTED ESTIMATION

*Muhammed O. Sayin, N. Denizcan Vanli, Suleyman S. Kozat*

Bilkent University, Ankara, Turkey

## ABSTRACT

We provide a logarithmic upper-bound on the regret function of the diffusion implementation for the distributed estimation. For certain learning rates, the bound shows guaranteed performance convergence of the distributed least mean square (DLMS) algorithms to the performance of the best estimation generated with hindsight of spatial and temporal data. We use a new cost definition for distributed estimation based on the widely-used statistical performance measures and the corresponding global regret function. Then, for certain learning rates, we provide an upper-bound on the global regret function without any statistical assumptions.

***Index Terms***— Regret, distributed, estimation, diffusion

## I. INTRODUCTION

Distributed network of nodes provides enriched observation ability over the monitored phenomena. In distributed estimation framework, we utilize this ability to estimate a parameter of interest by distributing the processing over the network. Diffusion implementation is one of the commonly used methods in distributed signal processing [1], [2]. Each node diffuses information to its neighbors and performs a local adaptive estimation algorithm more effectively with the benefit of the exchanged information [1]. In [1], nodes use the least mean square algorithm in local estimation and share the parameter estimate within a predefined neighborhood. The analysis of distributed estimation is rather challenging because of the cooperation of the nodes and in the literature authors provide performance analysis for certain statistical profiles [1], [2].

In this work, we avoid any statistical assumptions and aim to provide a deterministic performance analysis which is guaranteed to hold for any spatial or temporal data. To do that, we use a new cost definition for distributed estimation algorithms [3], which satisfies the global performance measures used in [1] and [2]. Each local parameter estimation is expected to converge to the optimum solution which yields the minimum cost for overall spatial and temporal data, i.e., the parameter of interest. Hence, the new cost also bills the performance of each local parameter estimate over the observations of any other nodes. Then, we use a global regret function, which is used as a performance measure in deterministic analysis excessively [4], [5]. We can define the regret of any algorithm as the difference between the cost of the algorithm and the minimum possible cost we could have with hindsight. Through the new cost and global regret definitions, we provide a logarithmic regret upper-bound on the performance of the diffusion based distributed estimation (specifically adapt-then-combine strategy [2]) for certain learning rates, which shows guaranteed performance for any spatial or temporal data.

## II. DIFFUSION IMPLEMENTATION

In a distributed network of $N$ nodes, each node $i$ observes a parameter of interest[1] $\mathbf{w_o} \in \mathbb{R}^p$ through a linear model

$$d_{i,t} = \mathbf{w_o}^T \mathbf{u}_{i,t} + v_{i,t},$$

where $i$ and $t$ are the node and time indices respectively. $v_{i,t}$ denotes the observation noise and $\mathbf{u}_{i,t} \in \mathbb{R}^p$ is the local regression vector.

In diffusion implementation framework, each node exchanges information with nodes from its neighborhood $\mathcal{N}_i$ and performs an estimation algorithm through the local observation $d_{i,t}$, the local regression vector $\mathbf{u}_{i,t}$ and the diffused information from the neighboring nodes. For example, the diffused information from $j$th node might be the local parameter estimation, i.e., $\mathbf{w}_{j,t}$, [1], [2]. In [2], authors examine the change of the performance of the algorithms with respect to the aggregation of the diffused information before and after the adaptation. They show that the adapt-then-combine (ATC) strategy outperforms the combine-then-adapt (CTA) strategy. Hence, in this paper, we provide the regret bound for the ATC strategy. The ATC update is given by

$$\phi_{i,t+1} = \mathbf{w}_{i,t} + \mu_i \mathbf{u}_{i,t} \left( d_{i,t} - \mathbf{u}_{i,t}^T \mathbf{w}_{i,t} \right)$$
$$\mathbf{w}_{i,t+1} = \sum_{j \in \mathcal{N}_i} \gamma_{i,j} \phi_{j,t+1}, \tag{1}$$

where $\mu_i > 0$ is the local step size and $\phi_{i,t+1}$ is an intermediate parameter vector. The combination weights for the parameter estimates are denoted by $\gamma_{i,j}$'s and the combination matrix $\mathbf{\Gamma}$ is given by

$$\mathbf{\Gamma} = \begin{bmatrix} \gamma_{11} & \cdots & \gamma_{1N} \\ \vdots & \ddots & \vdots \\ \gamma_{N1} & \cdots & \gamma_{NN} \end{bmatrix},$$

which is determined through certain combination rules, e.g., the Metropolis [6], with the constraint that $\mathbf{\Gamma 1} = \mathbf{1}$ for

---

[1] As notation we use bold lowercase (uppercase) letters for vectors (matrices). For a vector $\mathbf{u}$, $\mathbf{u}^T$ denotes its transpose and $\|\mathbf{u}\|$ is the $l$-2 norm.

unbiased convergence. In [1] and [2], the authors define the global performance measures for distributed estimation as follows:

$$\eta_t = \frac{1}{N} E \|\tilde{\mathbf{w}}_t\|^2, \tag{2}$$

$$\zeta_t = \frac{1}{N} E \|\mathbf{e}_{a,t}\|^2, \tag{3}$$

where $\tilde{\mathbf{w}}_t \triangleq \underline{\mathbf{w}}_{\mathbf{o}} - \underline{\mathbf{w}}_t$ is the global deviation vector, $\mathbf{e}_{a,t}$ is the global *a priori* error vector with the global parameters defined as

$$\begin{aligned}
\underline{\mathbf{w}}_{\mathbf{o}} &= \mathrm{col}\{\mathbf{w}_{\mathbf{o}}, \ldots, \mathbf{w}_{\mathbf{o}}\}_{Np \times 1} \\
\underline{\mathbf{w}}_t &= \mathrm{col}\{\mathbf{w}_{1,t}, \ldots, \mathbf{w}_{N,t}\}_{Np \times 1} \\
\mathbf{e}_{a,t} &= \mathrm{col}\{e_{a_{1,t}}, \ldots, e_{a_{N,t}}\}_{N \times 1}
\end{aligned} \tag{4}$$

and the local *a priori* error is $e_{a_{i,t}} = \mathbf{u}_{i,t}^T (\mathbf{w}_{\mathbf{o}} - \mathbf{w}_t)$. Note that (2) gives the global mean-square deviation and (3) yields the global excess mean square error.

In [1] and [2], authors provide performance analysis for distributed least squares algorithms under some assumptions for certain statistical profiles. In the following we provide a performance analysis for the diffusion implementation in the deterministic framework without any statistical assumption.

## III. LOGARITHMIC REGRET BOUND

With respect to the global performance measures (2) and (3), we expect the parameter estimations of all nodes to perform like $\mathbf{w}_*$, which is the best estimate we made if we would access to all spatial and temporal data overall network. Particularly, the estimation of each node should also perform well for the regression data of other nodes. Hence, the cost of the distributed estimation at time $T$ is given by

$$\mathrm{Cost}_T(\mathrm{DE}) = \frac{1}{N} \sum_{t=1}^{T} \sum_{i=1}^{N} \sum_{j=1}^{N} \left(d_{i,t} - \mathbf{u}_{i,t}^T \mathbf{w}_{j,t}\right)^2.$$

Note that in [3], authors use the same cost definition for the distributed autonomous online learning algorithm.

In the deterministic framework, regret function is a performance measure defined as the difference between the total cost and the cost of the best single decision, e.g., $\mathbf{w}_*$, which is chosen with the benefit of the hindsight [5]. We introduce a global regret function over the network as follows:

$$\begin{aligned}
\mathrm{Regret}_T(\mathrm{DE}) \triangleq & \frac{1}{N} \sum_{t=1}^{T} \sum_{j=1}^{N} \sum_{i=1}^{N} \left(d_{i,t} - \mathbf{u}_{i,t}^T \mathbf{w}_{j,t}\right)^2 \\
& - \sum_{t=1}^{T} \sum_{i=1}^{N} \left(d_{i,t} - \mathbf{u}_{i,t}^T \mathbf{w}_*\right)^2.
\end{aligned} \tag{5}$$

We define the cost function as

$$f_{i,t}(\mathbf{w}_{j,t}) \triangleq \left(d_{i,t} - \mathbf{u}_{i,t}^T \mathbf{w}_{j,t}\right)^2.$$

Then, (5) yields

$$\mathrm{Regret}_T(\mathrm{DE}) = \frac{1}{N} \sum_{t=1}^{T} \sum_{j=1}^{N} \sum_{i=1}^{N} \left[f_{i,t}(\mathbf{w}_{j,t}) - f_{i,t}(\mathbf{w}_*)\right].$$

We note that $f_{i,t}(\mathbf{w}_{j,t})$ is a convex function around $\mathbf{w}_{j,t}$, thus, the Hessian matrix $\nabla^2 f_{i,t}(\mathbf{w}_{j,t})$ is a positive semi-definite matrix, i.e., $\nabla^2 f_{i,t}(\mathbf{w}_{j,t}) \succeq \mathbf{0}$. The Hessian of the strictly convex cost functions is lower bounded by a number $H > 0$ if and only if

$$\nabla^2 f_{i,t}(\mathbf{w}_{j,t}) - H\mathbf{I}_p \succeq \mathbf{0}$$

is a positive semi-definite matrix. In [5], such functions are called *H-strong convex*. We can also upper bound the gradients of the cost function by a number $G$ provided that

$$\sup_{\mathbf{w} \in \mathbb{R}^p, t \in [T]} \|\nabla f_{i,t}(\mathbf{w}_{j,t})\| \leq G.$$

In addition, we assume that there are finite $A, D \in \mathbb{R}$ such that $\|\mathbf{u}_{i,t}\| < A$ and $|d_{i,t}| < D$ for all $i \in \{1, \cdots, N\}$ and $t$.

In [6], authors argue that the distributed linear averaging iterations converge to the average if and only if the combination matrix $\mathbf{\Gamma}$ yields

$$\lim_{t \to \infty} \mathbf{\Gamma}^t = \frac{\mathbf{1}\mathbf{1}^T}{N}.$$

This brings in the following constraints on $\mathbf{\Gamma}$: 1) $\mathbf{1}^T \mathbf{\Gamma} = \mathbf{1}^T$, 2) $\mathbf{\Gamma}\mathbf{1} = \mathbf{1}$, and 3) $\rho\left(\mathbf{\Gamma} - \frac{\mathbf{1}\mathbf{1}^T}{N}\right) < 1$, where $\rho(\cdot)$ denotes the spectral radius of the matrix. If the weights in $\mathbf{\Gamma}$ are non-negative, these conditions yields that $\mathbf{\Gamma}$ is doubly stochastic. Then, for aperiodic and irreducible $\mathbf{\Gamma}$, through the finite-state Markov chain theory, we have

$$\forall j \quad \sum_{i=1}^{N} \left| [\mathbf{\Gamma}^t]_{i,j} - \frac{1}{N} \right| \leq \theta \beta^t, \tag{6}$$

where $\theta > 0$ and $0 < \beta < 1$. In [3], authors set $\theta = 2$ and choose $\beta$ from the minimum nonzero values of $\mathbf{\Gamma}$.

We choose the same time dependent step size at all nodes and initialize each parameter estimate with the same value. Then, the following theorem provides a logarithmic bound on the regret function of ATC strategy for the doubly stochastic combination matrix $\mathbf{\Gamma}$.

**Theorem.** *The diffusion based distributed estimation with step sizes $\mu_{i,t+1} = \mu_{t+1} = \frac{1}{Ht}$ achieves the following guarantee, for all $T \geq 1$*

$$\mathrm{Regret}_T(\mathrm{DDE}) \leq \frac{G^2}{2H} C(1 + \log(T)), \tag{7}$$

*where*

$$C = N\left(1 + 2\frac{2G + AD}{G}\frac{\theta}{1 - \beta}\right).$$

In the next section, we provide the proof of the theorem.

## IV. PROOF OF THE THEOREM

The ATC strategy (1) leads the following updates:

$$\phi_{i,t+1} = \mathbf{w}_{i,t} - \mu_{t+1}\nabla f_{i,t}(\mathbf{w}_{i,t}), \tag{8}$$

$$\mathbf{w}_{i,t+1} = \sum_{j=1}^{N} \gamma_{i,j}\phi_{j,t+1}. \tag{9}$$

We can combine (8) and (9) as follows

$$\mathbf{w}_{i,t+1} = \sum_{j=1}^{N} \gamma_{i,j} \mathbf{w}_{j,t} - \mu_{t+1} \sum_{j=1}^{N} \gamma_{i,j} \nabla f_{j,t}(\mathbf{w}_{j,t}). \quad (10)$$

We assume that the combination matrix is doubly stochastic, i.e., $\sum_{i=1}^{N} \gamma_{i,j}$. Summing (10) from $i = 1$ to $N$, we obtain

$$\sum_{i=1}^{N} \mathbf{w}_{i,t+1} = \sum_{j=1}^{N} \mathbf{w}_{j,t} - \mu_{t+1} \sum_{j=1}^{N} \nabla f_{j,t}(\mathbf{w}_{j,t}). \quad (11)$$

We define an average parameter estimation vector $\bar{\mathbf{w}}_t$ as follows

$$\bar{\mathbf{w}}_t \triangleq \frac{1}{N} \sum_{i=1}^{N} \mathbf{w}_{i,t}.$$

Then, (11) yields

$$\bar{\mathbf{w}}_{t+1} = \bar{\mathbf{w}}_t - \mu_{t+1} \frac{1}{N} \sum_{i=1}^{N} \nabla f_{i,t}(\mathbf{w}_{i,t}). \quad (12)$$

Subtracting $\mathbf{w}_*$ from both side in (12) and taking $l_2$ norm square, we obtain

$$\sum_{i=1}^{N} \nabla f_{i,t}(\mathbf{w}_{i,t})^T (\bar{\mathbf{w}}_t - \mathbf{w}_*) \leq \frac{\mu_{t+1}}{2N} \left( \sum_{i=1}^{N} \|\nabla f_{i,t}(\mathbf{w}_{i,t})\| \right)^2$$
$$+ \frac{N}{2} \frac{\|\bar{\mathbf{w}}_t - \mathbf{w}_*\|^2 - \|\bar{\mathbf{w}}_{t+1} - \mathbf{w}_*\|^2}{\mu_{t+1}}, \quad (13)$$

where we use the triangular inequality as

$$\left\| \sum_{i=1}^{N} \nabla f_{i,t}(\mathbf{w}_{i,t}) \right\| \leq \sum_{i=1}^{N} \|\nabla f_{i,t}(\mathbf{w}_{i,t})\|.$$

The Taylor series expansion of the cost function $f_{i,t}(\cdot)$ leads

$$f_{i,t}(\bar{\mathbf{w}}_t) = f_{i,t}(\mathbf{w}_{j,t}) + \nabla f_{i,t}(\mathbf{w}_{j,t})^T (\bar{\mathbf{w}}_t - \mathbf{w}_{j,t})$$
$$+ \frac{1}{2}(\bar{\mathbf{w}}_t - \mathbf{w}_{j,t})^T \nabla^2 f_{i,t}(\mathbf{w}_{j,t})(\bar{\mathbf{w}}_t - \mathbf{w}_{j,t}) \quad (14)$$

and

$$f_{i,t}(\mathbf{w}_*) = f_{i,t}(\bar{\mathbf{w}}_t) + \nabla f_{i,t}(\bar{\mathbf{w}}_t)^T (\mathbf{w}_* - \bar{\mathbf{w}}_t)$$
$$+ \frac{1}{2}(\mathbf{w}_* - \bar{\mathbf{w}}_t)^T \nabla^2 f_{i,t}(\bar{\mathbf{w}}_t)(\mathbf{w}_* - \bar{\mathbf{w}}_t). \quad (15)$$

By (14) and (15), we get

$$\nabla f_{i,t}(\bar{\mathbf{w}}_t)^T (\bar{\mathbf{w}}_t - \mathbf{w}_*) \geq f_{i,t}(\mathbf{w}_{j,t}) - f_{i,t}(\mathbf{w}_*)$$
$$- \nabla f_{i,t}(\mathbf{w}_{j,t})^T (\mathbf{w}_{j,t} - \bar{\mathbf{w}}_t)$$
$$+ \frac{H}{2} \|\bar{\mathbf{w}}_t - \mathbf{w}_{j,t}\|^2 + \frac{H}{2} \|\bar{\mathbf{w}}_t - \mathbf{w}_*\|^2, \quad (16)$$

where the last two term on the right hand side (RHS) follows from the $H$-strong convexity.

We note that the term on the left hand side of (16) could be written as

$$\nabla f_{i,t}(\bar{\mathbf{w}}_t)^T (\bar{\mathbf{w}}_t - \mathbf{w}_*) = - \left[ \mathbf{u}_{i,t}(d_{i,t} - \mathbf{u}_{i,t}^T \bar{\mathbf{w}}_t) \right]^T$$
$$\times (\bar{\mathbf{w}}_t - \mathbf{w}_*)$$

and leads to

$$\nabla f_{i,t}(\mathbf{w}_{i,t})^T (\bar{\mathbf{w}}_t - \mathbf{w}_*) = \nabla f_{i,t}(\bar{\mathbf{w}}_t)^T (\bar{\mathbf{w}}_t - \mathbf{w}_*)$$
$$+ (\mathbf{w}_{i,t} - \bar{\mathbf{w}}_t)^T \mathbf{u}_{i,t} \mathbf{u}_{i,t}^T (\bar{\mathbf{w}}_t - \mathbf{w}_*) \quad (17)$$

Through (16), (17), and summing from $j = 1$ to $N$, we have

$$\nabla f_{i,t}(\mathbf{w}_{i,t})^T (\bar{\mathbf{w}}_t - \mathbf{w}_*) \geq \frac{1}{N} \sum_{j=1}^{N} [f_{i,t}(\mathbf{w}_j, t) - f_{i,t}(\mathbf{w}_*)]$$
$$+ \frac{H}{2N} \sum_{j=1}^{N} \|\mathbf{w}_{j,t} - \bar{\mathbf{w}}_t\|^2 + \frac{H}{2} \|\bar{\mathbf{w}}_t - \mathbf{w}_*\|^2$$
$$- \frac{1}{N} \sum_{j=1}^{N} \|\nabla f_{i,t}(\mathbf{w}_{j,t})\| \|\mathbf{w}_{j,t} - \bar{\mathbf{w}}_t\|$$
$$- \|\mathbf{u}_{i,t} \mathbf{u}_{i,t}^T (\bar{\mathbf{w}}_t - \mathbf{w}_*)\| \|\mathbf{w}_{i,t} - \bar{\mathbf{w}}_t\|. \quad (18)$$

We set a bound on the last term as

$$\|\mathbf{u}_{i,t} \mathbf{u}_{i,t}^T (\bar{\mathbf{w}}_t - \mathbf{w}_*)\| \leq \frac{1}{N} \sum_{j=1}^{N} (\|\nabla f_{i,t}(\mathbf{w}_{j,t})\| + A D)$$
$$\leq G + A D.$$

After some algebra, (13) and (18) yields

$$\frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{N} [f_{i,t}(\mathbf{w}_{j,t}) - f_{i,t}(\mathbf{w}_*)] \leq \frac{N}{2} \mu_{t+1} G^2$$
$$- \frac{H}{2} \sum_{i=1}^{N} \|\mathbf{w}_{i,t} - \bar{\mathbf{w}}_t\|^2 + (2G + A D) \sum_{i=1}^{N} \|\mathbf{w}_{i,t} - \bar{\mathbf{w}}_t\|$$
$$+ \frac{N}{2} \left[ \left( \frac{1}{\mu_{t+1}} - H \right) \|\bar{\mathbf{w}}_t - \mathbf{w}_*\|^2 - \frac{1}{\mu_{t+1}} \|\bar{\mathbf{w}}_{t+1} - \mathbf{w}_*\|^2 \right]. \quad (19)$$

In (19), we also have $\|\mathbf{w}_{i,t} - \bar{\mathbf{w}}_t\|$ terms. In [3], authors bound $\|\mathbf{w}_{i,t} - \bar{\mathbf{w}}_t\|$ terms using (6). The following lemma presents a similar result for the diffusion based distributed estimation.

**Lemma.** *For irreducible and aperiodic doubly stochastic combination matrix $\mathbf{\Gamma}$, the norm of the difference between the parameter estimate of any node, e.g., $\mathbf{w}_{i,t}$, and the average $\bar{\mathbf{w}}_t$ is bounded as follows:*

$$\|\mathbf{w}_{i,t} - \bar{\mathbf{w}}_t\| \leq NG\theta \sum_{\tau=1}^{t-1} \mu_{t-\tau+1} \beta^\tau.$$

**Proof.** We resort to the global parameter estimation $\underline{\mathbf{w}}_t$ definition (4) and define

$$\underline{\mathbf{f}}_t \triangleq \text{col}\{\nabla f_{1,t}(\mathbf{w}_{1,t}), \cdots, \nabla f_{N,t}(\mathbf{w}_{N,t})\}.$$

Then, by (10), we obtain

$$\underline{\mathbf{w}}_{t+1} = \underline{\boldsymbol{\Gamma}}\underline{\mathbf{w}}_t - \mu_{t+1}\underline{\boldsymbol{\Gamma}}\underline{\mathbf{f}}_t, \tag{20}$$

where $\underline{\boldsymbol{\Gamma}} = \boldsymbol{\Gamma} \otimes \mathbf{I}_p$. The iteration of (20) leads

$$\underline{\mathbf{w}}_t = \underline{\boldsymbol{\Gamma}}^{t-1}\underline{\mathbf{w}}_1 - \sum_{\tau=1}^{t-1} \mu_{t-\tau+1}\underline{\boldsymbol{\Gamma}}^{\tau}\underline{\mathbf{f}}_{t-\tau}. \tag{21}$$

We introduce $\mathbf{e} \triangleq \mathrm{col}\{1, \cdots, 1\} \otimes \mathbf{I}_p$ and $\mathbf{e}_i \triangleq \mathrm{col}\{0, \cdots, 1, \cdots, 0\} \otimes \mathbf{I}_p$ where only the $i$th term is 1. Since $\boldsymbol{\Gamma}$ is a right-stochastic matrix, i.e., $\mathbf{e}\boldsymbol{\Gamma} = \mathbf{e}$, through (21), we can bound the term $\|\bar{\mathbf{w}}_t - \mathbf{w}_{i,t}\|$ as follows

$$\|\bar{\mathbf{w}}_t - \mathbf{w}_{i,t}\| = \left\| \left(\frac{1}{N}\mathbf{e} - \mathbf{e}_i\right)\underline{\mathbf{w}}_t \right\|$$

$$\leq \left\| \left(\frac{\mathbf{e}}{N} - \mathbf{e}_i\right)\underline{\boldsymbol{\Gamma}}^{t-1}\underline{\mathbf{w}}_1 \right\| + \sum_{\tau=1}^{t-1} \mu_{t-\tau+1}\left\| \left(\frac{\mathbf{e}}{N} - \mathbf{e}_i\right)\underline{\boldsymbol{\Gamma}}^{\tau}\underline{\mathbf{f}}_{t-\tau} \right\|$$

$$\leq \|\bar{\mathbf{w}}_1 - \mathbf{w}_{i,1}\| + \sum_{\tau=1}^{t-1} \mu_{t-\tau+1}\|\underline{\mathbf{f}}_{t-\tau}\|\left\| \left(\frac{1}{N}\mathbf{e} - \mathbf{e}_i\right)\underline{\boldsymbol{\Gamma}}^{\tau} \right\|.$$

We assume that all parameter estimation vectors are initialized with the same value, i.e., $\bar{\mathbf{w}}_1 = \sum_{i=1}^{N} \mathbf{w}_{i,1} = \mathbf{w}_{i,1}$, then the difference term $\|\bar{\mathbf{w}}_1 - \mathbf{w}_{i,1}\|$ goes to zero. We also note that

$$\|\underline{\mathbf{f}}_{t-\tau}\| = \sum_{i=1}^{N} \|\nabla f_{i,t}(\mathbf{w}_{i,t})\| \leq NG.$$

Finally, by (6), we have

$$\left\| \frac{1}{N}\mathbf{e}\underline{\boldsymbol{\Gamma}}^{\tau} - \mathbf{e}_i\underline{\boldsymbol{\Gamma}}^{\tau} \right\| = \sum_{j=1}^{N} \left| [\underline{\boldsymbol{\Gamma}}^{\tau}]_{j,i} - \frac{1}{N} \right| \leq \theta\beta^{\tau}.$$

The proof is concluded. $\qquad\square$

Through the Lemma, the summation of (19) from $t = 1$ to $T$ leads

$$\frac{1}{N} \sum_{t=1}^{T} \sum_{i=1}^{N} \sum_{j=1}^{N} [f_{i,t}(\mathbf{w}_{j,t}) - f_{i,t}(\mathbf{w}_*)] \leq \frac{NG^2}{2} \sum_{t=1}^{T} \mu_{t+1}$$

$$+ \frac{N}{2} \sum_{t=1}^{T} \left[ \left(\frac{1}{\mu_{t+1}} - H\right)\|\bar{\mathbf{w}}_t - \mathbf{w}_*\|^2 - \frac{1}{\mu_{t+1}}\|\bar{\mathbf{w}}_{t+1} - \mathbf{w}_*\|^2 \right]$$

$$+ NG\theta(2G + AD) \sum_{t=1}^{T} \sum_{\tau=1}^{t-1} \mu_{t-\tau+1}\beta^{\tau}. \tag{22}$$

We drop $\|\mathbf{w}_{i,t} - \bar{\mathbf{w}}_t\|^2$ term in (22). This expands the upper bound on the regret function, however, results in simpler

bound expression. The last term on the RHS of (22) yields

$$\sum_{t=1}^{T} \left[ \left(\frac{1}{\mu_{t+1}} - H\right)\|\bar{\mathbf{w}}_t - \mathbf{w}_*\|^2 - \frac{1}{\mu_{t+1}}\|\bar{\mathbf{w}}_{t+1} - \mathbf{w}_*\|^2 \right]$$

$$= \left(\frac{1}{\mu_2} - H\right)\|\bar{\mathbf{w}}_1 - \mathbf{w}_*\|^2 - \underbrace{\frac{1}{\mu_2}\|\bar{\mathbf{w}}_2 - \mathbf{w}_*\|^2}$$

$$+ \underbrace{\left(\frac{1}{\mu_3} - H\right)\|\bar{\mathbf{w}}_2 - \mathbf{w}_*\|^2} - \frac{1}{\mu_3}\|\bar{\mathbf{w}}_3 - \mathbf{w}_*\|^2$$

$$\vdots$$

$$+ \left(\frac{1}{\mu_{T+1}} - H\right)\|\bar{\mathbf{w}}_T - \mathbf{w}_*\|^2 - \frac{1}{\mu_{T+1}}\|\bar{\mathbf{w}}_{T+1} - \mathbf{w}_*\|^2 \tag{23}$$

Re-arranging the sum such that the terms with the same time indices gathered together, we obtain

$$\sum_{t=1}^{T} \left[ \left(\frac{1}{\mu_{t+1}} - H\right)\|\bar{\mathbf{w}}_t - \mathbf{w}_*\|^2 - \frac{1}{\mu_t}\|\bar{\mathbf{w}}_t - \mathbf{w}_*\|^2 \right]. \tag{24}$$

Note that during the rearrangement of the sum we set $\frac{1}{\mu_1} = 0$ ($\mu_1$ is not used in the update (10)) and extend the upper-bound by neglecting the last term in (23). (24) implies that for $\mu_{t+1} = \frac{1}{Ht}$, the second term on the RHS of (22) goes to zero.

In [3], authors show that

$$\sum_{t=1}^{T} \sum_{\tau=1}^{t-1} \mu_{t-\tau+1}\beta^{\tau} \leq \frac{1}{1-\beta} \sum_{t=1}^{T} \mu_{t+1}.$$

Thus, for $\mu_{t+1} = \frac{1}{Ht}$, we have

$$\mathrm{Regret}_T(\mathrm{DDE}) \leq \left(\frac{NG^2}{2} + \frac{NG\theta(2G + AD)}{1-\beta}\right) \sum_{t=1}^{T} \frac{1}{Ht}$$

and $\sum_{t=1}^{T} \frac{1}{t} \leq 1 + \log(T)$. This completes the proof of the Theorem (7).

## V. CONCLUDING REMARKS

Diffusion implementation has appealed interest in the distributed estimation and provides improved convergence performance over the non-coherent update. In this paper, we provide a logarithmic regret upper bound on the diffusion based distributed estimation algorithms for certain learning rates. An upper bound on regret function is of interest because averaging the regret over time, we observe that logarithmic upper-bound goes to zero. This implies that the performance of the distributed estimation asymptotically converges to the performance of the best solution we could get with the hindsight of all spatial and temporal data.

## VI. REFERENCES

[1] C. G. Lopes and A. H. Sayed, "Diffusion least-mean squares over adaptive networks: Formulation and performance analysis," *IEEE Transactions on Signal Processing*, vol. 56, no. 7, pp. 3122–3136, 2008.

[2] F. S. Cattivelli and A. H. Sayed, "Diffusion LMS strategies for distributed estimation," *IEEE Transactions on Signal Processing*, vol. 58, no. 3, pp. 1035–1048, 2010.

[3] 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, no. 11, pp. 2483–2493, 2013.

[4] M. Zinkevich, "Online convex programming and generalized infinitesimal gradient ascent," in *Proceedings of the Twentieth International Conference (ICML)*, 2003, pp. 928–936.

[5] E. Hazan, A. Agarwal, and S. Kale, "Logarithmic regret algorithms for online convex optimization," *Mach. Learn.*, vol. 69, no. 2-3, pp. 169–192, Dec. 2007.

[6] Lin Xiao and S. Boyd, "Fast linear iterations for distributed averaging," *Syst. Control Lett.*, vol. 53, no. 1, pp. 65–78, 2004.