

# Exponentially Fast Parameter Estimation in Networks Using Distributed Dual Averaging<sup>†</sup>

Shahin Shahrampour<sup>‡</sup> and Ali Jadbabaie<sup>‡</sup>

**Abstract**—In this paper we present an optimization-based view of distributed parameter estimation and observational social learning in networks. Agents receive a sequence of random, independent and identically distributed (i.i.d.) signals, each of which individually may not be informative about the underlying true state, but the signals together are globally informative enough to make the true state identifiable. Using an optimization-based characterization of Bayesian learning as proximal stochastic gradient descent (with Kullback-Leibler divergence from a prior as a proximal function), we show how to efficiently use a distributed, online variant of Nesterov's dual averaging method to solve the estimation with purely local information. When the true state is globally identifiable, and the network is connected, we prove that agents eventually learn the true parameter using a randomized gossip scheme. We demonstrate that with high probability the convergence is exponentially fast with a rate dependent on the KL divergence of observations under the true state from observations under the second likeliest state. Furthermore, our work also highlights the possibility of learning under continuous adaptation of network which is a consequence of employing constant, unit stepsize for the algorithm.

## I. INTRODUCTION

Distributed estimation, detection, and observational social learning has been an intense focus of research over the past 3 decades [1]–[9], with applications ranging from sensor networks to social and economic networks. In these scenarios, agents in a network need to learn the value of a parameter, that might represent a state or decision (often called the state of the world), but each individual agent lacks the necessary information to estimate the state on its own. Instead, the global spread of information in the network provides agents with adequate data for recovering the true state and as a result, agents iteratively exchange information with their neighbors. In distributed sensor and robotic networks, agents use local diffusion to augment their imperfect observations with information from their neighbors [6], [10]–[14].

On the other hand, recent developments in distributed optimization have led to many advances and interesting decentralized algorithms, generalizing these results and at the same time opening new venues for development of principled distributed estimation algorithms. Examples of such papers include the works of researchers such as Nedić and Ozdaglar [15], Lobel and Ozdaglar [16], Ram, Nedić and Veeravalli [17] and Nedić, Olshevsky, Ozdaglar and

Tsitsiklis [18], Lopes and Sayed [19], and more recently the results of Duchi, Agarwal and Wainwright [20]. Of particular importance to the work in this paper is the work of Duchi *et al.* in [20] where the authors develop a distributed method based on dual averaging of subgradients. Using proper diminishing stepsize rule, their algorithm converges to the optimal solution in deterministic network change as well as stochastic.

The goal of this paper is to provide an optimization-based formulation of parameter estimation and social learning and develop a link between the two. Our motivation for the current study is the recent results of [7] and [8] in which the authors develop non-Bayesian learning schemes to circumvent the complexities associated with fully Bayesian estimation [1], as well as the results of [20] on distributed optimization. The proposed algorithms in [7] and [8], involve agents that repeatedly receive heterogeneous, private, random i.i.d. signals generated from a global likelihood function and the goal of agents is to learn the true state of the world using local marginals. Both papers show that under mild assumptions all agents eventually estimate the true parameter correctly. In [8], agents update their prior beliefs using private observations and then compute a weighted average of their beliefs with that of their neighbors, while in [7], agents update the logarithms of their beliefs using the local log-likelihood function. In both cases, under mild assumptions agents eventually learn the true state. We show that the results of [7], have a very interesting optimization-based interpretation. Exploring this connection and building an optimization-based rationale helps us quantify the pros and cons of different approaches to the problem.

A key unifying observation that links both recursive Bayesian learning and Maximum Likelihood Estimation (MLE) problems to online optimization (even in the centralized setting) is the view of MLE in Bayesian framework as an optimization in which the inner product of the belief vector and the global log-likelihood function (represented as a vector) is maximized, subject to the belief vector being a probability distribution over the space of parameters. Perhaps less well-known, is the fact that Bayesian parameter estimation can be derived from the exact same setup if the Kullback-Leibler divergence from a prior belief is added to the optimization cost or used as a *proximal function*. We show an efficient distributed counterpart of this idea using a stochastic variant of Nesterov's projected dual averaging [21]. Aggregating their private log-likelihood functions, agents average their local information, and in the same time step update estimates of the centralized beliefs in a step akin

<sup>†</sup>This work was supported by AFOSR MURI CHASE, and ONR BRC Program on Decentralized, Online Optimization.

<sup>‡</sup>The authors are with the Department of Electrical and Systems Engineering and General Robotics, Automation, Sensing and Perception (GRASP) Laboratory, University of Pennsylvania, Philadelphia, PA 19104-6228 USA. {shahin, jadbabai}@seas.upenn.edu

to applying Bayes rule on the aggregated log-likelihoods. When the true state is globally identifiable and the network is connected, we show that agents reach consensus on the beliefs in probability. More specifically, we prove that with high probability the convergence is exponentially fast with a rate dependent on the *average expected discrimination information* for the true state over the second likeliest state captured by the KL divergence of the observations under two aforementioned states. We further show that indeed there is no need for a diminishing stepsize rule as in general subgradient approaches, and a fixed stepsize of 1 can be used. Interestingly, the method recovers the distributed MAP algorithm proposed by [7] as a special case.

The rest of the paper is organized as follows. In the next section, we introduce the model under which agents interact, define our learning problem and formulate it as a constrained maximization. In section III, we recover Bayesian estimation with dual averaging. In section IV we show applying gossip distributed dual averaging under constant, unit stepsize rules results in learning in the probability sense, and the convergence is exponentially fast. Section V concludes.

## II. PRELIMINARIES

### A. Agents and Observation

We consider a network consisting of a finite number of agents  $V = \{1, 2, \dots, n\}$ . The agents indexed by  $i \in V$  seek a fixed, unique, true state of the world  $\theta^* \in \Theta$  with  $\Theta = \{\theta_1, \theta_2, \dots, \theta_m\}$  denoting a finite set of possible states. At each time  $t \geq 0$ , belief of agent  $i$  is denoted by  $\mu_{i,t}(\theta) \in \Delta\Theta$ , where  $\Delta\Theta$  is a probability distribution over the set  $\Theta$ . In particular,  $\mu_{i,0}(\theta) \in \Delta\Theta$  denotes the prior belief of agent  $i$  about the states of the world. For each agent  $i$ , we assume the prior  $\mu_{i,0}$  is in the interior of the probability simplex and as a result has no zero elements<sup>1</sup>.

The learning model is given by a conditional likelihood function  $\ell(s^t|\theta_j)$  which is governed by a state of the world  $\theta_j \in \Theta$ . The signal  $s_t = (s_1^t, s_2^t, \dots, s_n^t) \in S_1 \times \dots \times S_n$  is generated at each time  $t$ , and  $s_i^t \in S_i$  denotes the signal privately observed by agent  $i$  at time  $t$ , where  $S_i$  is the signal space for agent  $i$ .  $\ell_i(\cdot|\theta_j)$  represents the  $i$ -th marginal of  $\ell(\cdot|\theta_j)$ , and we let the vector  $\ell_i(\cdot|\theta) = [\ell_i(\cdot|\theta_1), \dots, \ell_i(\cdot|\theta_m)]^T$ , for any  $i \in V$ , where  $\ell_i(\cdot|\theta_j) > 0$  for all signals at all times. Agent  $i$  at time  $t$ , has access to the parametrized likelihood of the realized private signal  $s_i^t$ , i.e., it knows the value of  $\ell_i(s_i^t|\theta)$ , but does not have access to the likelihood functions of other agents, i.e., it does not know  $\ell_j(\cdot|\theta)$  for any  $j \neq i$ . Generated signals are i.i.d. over time and also independent over agents. We also define  $\bar{\Theta}_i$  as the set of states that are observationally equivalent to  $\theta^*$  for agent  $i$ ; in other words,  $\bar{\Theta}_i = \{\theta_j \in \Theta : \ell_i(s_i|\theta_j) = \ell_i(s_i|\theta^*) \forall s_i \in S_i\}$  with probability one. Let  $\bar{\Theta} = \bigcap_{i=1}^n \bar{\Theta}_i$  be the set of states that are observationally equivalent to  $\theta^*$  from all agents perspective. We assume

- A1. The true state is globally identifiable, and hence,  $\bar{\Theta} = \{\theta^*\}$ .
- A2. Each log-marginal  $\log \ell_i(\cdot|\theta_j)$  has a bounded variance.

The probability triple  $(\Omega, \mathcal{F}, \mathbb{P}^{\theta^*})$  is defined such that  $\Omega = (\otimes_{i=1}^n S_i)^{\mathbb{N}}$ ,  $\mathcal{F}_t$  is the smallest  $\sigma$ -field containing the information about all agents up to time  $t$ , and  $\mathbb{P}^{\theta^*}$  is the true probability measure with respect to  $\Omega$  with  $\mathbb{E}^*$  being its corresponding expectation operator.  $\mathbb{N}$  represents the natural numbers and  $\mathcal{F} = \bigcup_{t=1}^{\infty} \mathcal{F}_t$ .

*Definition 1:* Agent  $i \in V$  asymptotically learns the true parameter  $\theta^*$  on a path  $\{s^t\}_{t=1}^{\infty}$  if, along that path,

$$\mu_{i,t}(\theta^*) \rightarrow 1 \quad \text{as } t \rightarrow \infty.$$

The definition is intuitive as learning occurs when agents assign probability one to the unique true parameter.

### B. Time Model and Communication Structure

The interaction between agents is captured by an undirected graph  $G = (V, E)$ , where  $V$  is the set of agents and if there is a link between agent  $i$  and agent  $j$ , the pair  $\{i, j\}$  belongs to the set  $E$ . We let  $\mathcal{N}_i = \{j \in V : \{i, j\} \in E\}$  be the set of neighbors of agent  $i$ .

Agents communication conforms to an invariant gossip algorithm [22], wherein each node has a clock which ticks according to a rate 1 Poisson process. Equivalently, there is a single global clock which ticks according to a rate  $n$  Poisson process at times  $T_t$ , where  $\{T_t - T_{t-1}\}$  are i.i.d. exponential random variables with rate  $n$ . In the analysis, we use the index  $t$  to refer to the  $t$ -th time slot  $[T_{t-1}, T_t)$ ,  $t \geq 0$ . At each tick  $T_t$  of the global clock, agent  $I_t \in V$  is picked uniformly at random. Then, it contacts a neighbor  $J_t \in V$  with probability  $P_{I_t J_t}$ , and they update their belief. Denoting the communication matrix by  $W(t)$ , this amounts to  $W(t)$  taking the form

$$W(t) = I - \frac{(\mathbf{e}_{I_t} - \mathbf{e}_{J_t})(\mathbf{e}_{I_t} - \mathbf{e}_{J_t})^T}{2},$$

with probability  $\frac{1}{n}P_{I_t J_t}$ , where  $\mathbf{e}_i$  is the  $i$ -th unit vector in the standard basis of  $\mathbb{R}^n$ . Hence, the matrix  $P = [P_{ij}]$  has nonnegative entries and  $P_{ij} > 0$  only if  $\{i, j\} \in E$ . By definition,  $P$  is row stochastic with largest eigenvalue 1. We assume

- A3. The network is *connected* i.e., there exists a path from any agent  $i$  to any agent  $j$ , and the second largest eigenvalue of  $\mathbb{E}[W(t)]$  is strictly less than one in magnitude.

The connectivity constraint in assumption (A3) guarantees the information flow in the network. The assumption, for instance, holds if the underlying structure of the network is connected and nonbipartite.

### C. Problem Setup and Formulation

The MLE problem of finding the likeliest true state, can be formulated in terms of a belief vector  $\mu$  as the following

<sup>1</sup>We will see that this assumption is just for dealing with log-likelihood functions and technical issues; otherwise, we only need strict positivity of beliefs over the true state.

optimization

$$\max_{\mu \in \Delta\Theta} \left\{ f(\mu) \triangleq \mu^T \sum_{i=1}^n \mathbb{E}_{s_i}^* [\log \ell_i(s_i|\theta)] \right\}, \quad (1)$$

with  $\mu^*$  being its optimal solution. In the next section we discuss that a regularization term can be added to the objective function of (1) or used as a common *proximal function* among the agents. Alternatively, one might cast (1) as a quest for finding the MLE solution

$$\theta^* = \operatorname{argmax}_{\theta_j \in \Theta} \left\{ \mathbb{E}_s^* [\log \ell(s|\theta_j)] \right\}. \quad (2)$$

The equivalence of (1) and (2) follows immediately from the independence of agents observations, and the global identifiability of  $\theta^*$  (assumption A1) which guarantees that (1) has a unique maximizer. In the sequel, without loss of generality, we assume the components of the vector  $\mathbb{E}_s^* \log \ell(s|\theta)$  are in descending order, i.e.

$$\mathbb{E}_s^* [\log \ell(s|\theta_1)] > \mathbb{E}_s^* [\log \ell(s|\theta_2)] \geq \dots \geq \mathbb{E}_s^* [\log \ell(s|\theta_m)], \quad (3)$$

where the strict inequality on the left-hand side of (3) is due to uniqueness of  $\theta^* = \theta_1$ . Hence,  $\theta_1$  is the unique true state that is aimed to be recovered, and  $\mu^* = e_1$ .

### III. BAYESIAN ESTIMATION VIA NESTEROV'S DUAL AVERAGING

The learning problem formulated in (1) is a maximization over a closed convex set, so the structure of the problem allows us to apply a distributed generalization of the centralized dual averaging method proposed in [21]. First, however, we show how Bayesian learning can be viewed with an optimization lens.

A common approach to tackle problem (1) is to consider the empirical average as the cost function, and solve the online stochastic learning problem. To this end, we employ a regularized dual averaging scheme generating a sequence of iterates  $\{\mu_t, z_t\}_{t=0}^\infty$ , where  $\mu_t \in \Delta\Theta$  and  $z_t \in \mathbb{R}^m$ . At time period  $t$  the algorithm receives  $g_t$ , the stochastic gradient of the objective function, and performs the following set of centralized updates:

$$z_{t+1} = z_t + g_t \quad \text{and} \quad \mu_{t+1} = \prod_{\Delta\Theta}^{\psi}(z_{t+1}, \alpha_t), \quad (4)$$

where  $\{\alpha_t\}_{t=0}^\infty$  is a non-increasing sequence of positive stepsize,  $\psi(\cdot)$  is a so called *proximal function*, and

$$\prod_{\Delta\Theta}^{\psi}(z, \alpha) \triangleq \operatorname{argmin}_{x \in \Delta\Theta} \left\{ -\langle z, x \rangle + \frac{1}{\alpha} \psi(x) \right\}, \quad (5)$$

with  $\langle z, x \rangle$  being the standard inner product in the space of  $\mathbb{R}^m$ .

The *dual* update  $z$ , essentially integrates the stochastic gradients, and the second update projects the integration on the feasible set while regularizing the projection using a proximal function.

A particularly relevant example of a proximal function is the Kullback-Leibler (KL) divergence (also known as relative entropy) from an initial belief  $\mu_0$  defined as [23]:

$$\psi(x) = D_{KL}(x||\mu_0) \triangleq \sum_{i=1}^m [x]_i \log \frac{[x]_i}{\mu_0(\theta_i)}, \quad (6)$$

for any  $x \in \Delta\Theta$ , where  $[x]_i$  is the  $i$ -th component of the vector  $x$ . It is straightforward to verify that KL divergence from  $\mu_0$  is strongly convex with respect to the  $\ell_1$ -norm on the probability simplex  $\{x|x \geq 0, \sum_{i=1}^m [x]_i = 1\}^2$ .

The following proposition shows how the set of updates (4) equipped with the KL divergence could be viewed as an optimization counterpart of Bayesian rule.

*Proposition 2:* Given update rules (4) with stepsize sequence  $\{\alpha_t = 1\}_{t=0}^\infty$ , using KL divergence as the proximal function, following the stochastic gradient at each time period  $t$ , and letting  $z_0 = 0$ , we obtain the Bayes rule as

$$\mu_t(\theta) = \frac{\mu_{t-1}(\theta) \odot \ell(s^t|\theta)}{\sum_{j=1}^m \mu_{t-1}(\theta_j) \ell(s^t|\theta_j)}, \quad (7)$$

where  $\odot$  is component-wise multiplication.

*Proof:* To solve (1) with updates (4), since the stochastic gradient is  $g_t = \sum_{i=1}^n \log \ell_i(s_i^t|\theta)$ , performing the first update, we have

$$z_t = \sum_{i=1}^n \sum_{\tau=0}^{t-1} \log \ell_i(s_i^\tau|\theta) = \sum_{\tau=0}^{t-1} \log \ell(s^\tau|\theta).$$

Using  $\psi(x) = \sum_{j=1}^m [x]_j \log \frac{[x]_j}{\mu_0(\theta_j)}$  as the proximal function, we need to solve

$$\mu_t = \operatorname{argmin}_{x \in \Delta\Theta} \left\{ -x^T z_t + \sum_{j=1}^m [x]_j \log \frac{[x]_j}{\mu_0(\theta_j)} \right\}. \quad (8)$$

Leaving the positivity constraint implicit, we can write (8) as the maximization of the following Lagrangian

$$\mathbb{L}(x, \lambda) = x^T \sum_{\tau=0}^{t-1} \log \ell(s^\tau|\theta) - \sum_{j=1}^m [x]_j \log \frac{[x]_j}{\mu_0(\theta_j)} + \lambda(x^T \mathbf{1} - 1), \quad (9)$$

where  $\mathbf{1}$  is the vector of all ones. Differentiating (9) we get

$$\begin{aligned} \frac{\partial}{\partial [x]_j} \mathbb{L}(x, \lambda) &= \sum_{\tau=0}^{t-1} \log \ell(s^\tau|\theta_j) - \log [x]_j + \log \mu_0(\theta_j) - 1 + \lambda \\ \frac{\partial}{\partial \lambda} \mathbb{L}(x, \lambda) &= x^T \mathbf{1} - 1. \end{aligned}$$

Setting the above equations to zero, we get

$$x = \exp^{\lambda-1} \mu_0 \odot \prod_{\tau=0}^{t-1} \ell(s^\tau|\theta) \quad (10)$$

$$x^T \mathbf{1} = 1, \quad (11)$$

and replacing  $x$  in (11) by (10) we have

$$\exp^{\lambda-1} = \frac{1}{\sum_{j=1}^m \mu_0(\theta_j) \prod_{\tau=0}^{t-1} \ell(s^\tau|\theta_j)}. \quad (12)$$

<sup>2</sup>At origin we consider the limit; in other words, we define  $0 \log(0) = 0$ .

Hence, by (10) and (12) we have

$$\mu_t(\theta) = \frac{\mu_0(\theta) \odot \prod_{\tau=0}^{t-1} \ell(s^\tau | \theta)}{\sum_{j=1}^m \mu_0(\theta_j) \prod_{\tau=0}^{t-1} \ell(s^\tau | \theta_j)},$$

and (7) follows by compositionality of the above equation. ■

We derived a closed-form solution for  $\mu_t(\theta)$  that essentially performs the Bayesian update; each agent aggregates information up to time  $t$ , and then, infers the posterior from prior. One can prove the almost sure convergence of  $\mu_t(\theta)$  combining the arguments in [24] and [8]. However, we are interested in solving (1) in a decentralized manner, and we only use a generalized result of Proposition 2 later.

#### IV. DISTRIBUTED STOCHASTIC LEARNING

We now show that the centralized optimization studied in the previous section can be distributed over the network. Contrary to the centralized algorithm, each agent  $i \in V$ , at  $t$ -th slot only observes  $g_{i,t}$ , the stochastic gradient of its associated log-likelihood function, while it does not have access to the signals of other agents. The communication structure is based on a randomized gossip scheme. Let the global Poisson clock at the beginning of the  $t$ -th slot tick for agent  $i$  (with probability  $\frac{1}{n}$ ), and let agent  $i$  contact a neighboring node  $j$  (with probability  $P_{ij}$ ). Then, agents  $i$  and  $j$  average their accumulated observations from previous slots, and add their new stochastic gradients to form the following online updates

$$z_{i,t+1} = \frac{z_{i,t} + z_{j,t}}{2} + g_{i,t} \quad \text{and} \quad z_{j,t+1} = \frac{z_{i,t} + z_{j,t}}{2} + g_{j,t}, \quad (13)$$

while in that slot, any other agent  $k \notin \{i, j\}$  does not contact its neighbors, and only follows its own stochastic gradient  $g_{k,t}$ , so we have

$$z_{k,t+1} = z_{k,t} + g_{k,t}. \quad (14)$$

Having updated their observations, all agents calculate their estimates

$$\mu_{i,t+1}(\theta) = \prod_{\Delta \in \Theta}^{\psi} (z_{i,t+1}, \alpha_t), \quad (15)$$

where  $\prod_{\Delta \in \Theta}^{\psi}(z, \alpha)$  is previously defined in (5). Letting

$$Z_t = \begin{pmatrix} z_{1,t} \\ z_{2,t} \\ \vdots \\ z_{n,t} \end{pmatrix} \quad \text{and} \quad G_t = \begin{pmatrix} g_{1,t} \\ g_{2,t} \\ \vdots \\ g_{n,t} \end{pmatrix},$$

the set of updates (13) and (14) can be represented in the matrix form as follows

$$Z_{t+1} = \tilde{W}(t) Z_t + G_t, \quad (16)$$

where  $\tilde{W}(t) = W(t) \otimes \mathbf{I}_{m \times m}$ , and the random matrix  $W(t)$  with probability  $\frac{1}{n} P_{ij}$  takes the form

$$W(t) = I - \frac{(\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^T}{2}. \quad (17)$$

We use the above distributed stochastic scheme to optimize (1) equipped with the KL divergence defined in (6). It is noteworthy that in the distributed setting, employing the KL divergence from the initial belief, each agent exhibits inertia to a default opinion over the states. We prove in the next lemma that  $\mu_{i,t}(\theta)$  preserves a Bayes-like evolution.

*Lemma 3:* Given the set of update rules (13)-(14)-(15) with stepsize sequence  $\{\alpha_t = 1\}_{t=0}^{\infty}$ , following its stochastic gradient  $g_{i,t} = \log \ell_i(s_i^t | \theta)$  at  $t$ -th time period, if we let  $z_{i,0} = 0^3$ , agent  $i$ 's estimator evolves as

$$\mu_{i,t}(\theta) = \frac{\mu_{i,0}(\theta) \odot \exp[t\Phi_{i,t}(\theta)]}{\sum_{j=1}^m \mu_{i,0}(\theta_j) \exp[t\Phi_{i,t}(\theta_j)]}, \quad (18)$$

where

$$\Phi_{i,t}(\theta) = \frac{1}{t} \sum_{\tau=0}^{t-1} \sum_{k=1}^n \left[ \prod_{\rho=1}^{t-1-\tau} W(t-\rho) \right]_{ik} \log \ell_k(s_k^\tau | \theta). \quad (19)$$

*Proof:* The discrete-time linear system (16) has the closed-form solution

$$Z_t = \left( \prod_{\rho=1}^t \tilde{W}(t-\rho) \right) Z_0 + \sum_{\tau=0}^{t-1} \left( \prod_{\rho=1}^{t-\tau-1} \tilde{W}(t-\rho) \right) G_\tau.$$

Letting  $Z_0 = 0$ , since  $\tilde{W}(t) = W(t) \otimes \mathbf{I}_{m \times m}$ , by basic properties of Kronecker product, we can extract  $z_{i,t}$  from  $Z_t$  for each  $i$  to get

$$z_{i,t} = \sum_{\tau=0}^{t-1} \sum_{k=1}^n \left[ \prod_{\rho=1}^{t-1-\tau} W(t-\rho) \right]_{ik} \log \ell_k(s_k^\tau | \theta) = t\Phi_{i,t}(\theta).$$

We now need to solve (15) to complete the proof. The argument follows in the same fashion as Proposition 2. Forming the Lagrangian as in (9), writing the first order conditions using  $z_{i,t}$  derived above, and following the same steps as in (10), (11) and (12), the closed-form solution (18) follows immediately. ■

Equation (18) shows that at each time period  $t \geq 0$ , the set of distributed update rules (13)-(14)-(15) construct a *Gibbs distribution* over the states. As we shall see in the next subsection,  $\Phi_{i,t}(\theta)$  plays a key role on the convergence of the sequence of distributions generated over time.

#### A. Convergence Analysis

We now exhibit that aggregating information over time, agents have arbitrarily close opinions in a connected network. This is captured by the fact that the limit of (19) is independent of  $i$  which indexes agents. In this direction, we study the limit behavior of (19) in the following lemma.

*Lemma 4:* Under assumptions (A2) and (A3), the vector  $\Phi_{i,t}(\theta)$  defined in (19) converges in the probability sense as follows

$$\Phi_{i,t}(\theta) \xrightarrow{p} \Phi_\infty(\theta) = \frac{1}{n} \sum_{k=1}^n \mathbb{E}_{s_k}^* [\log \ell_k(s_k | \theta)].$$

*Proof:* The sequence  $\{W(t)\}_{t=0}^{\infty}$  is doubly stochastic, and the product term in  $\Phi_{i,t}(\theta)$  preserves doubly stochasticity, so for any  $j \in \{1, 2, \dots, m\}$  we have

<sup>3</sup>Regardless of the zero initial value, all the results hold asymptotically, and this condition only simplifies our derivation.

$$\begin{aligned}\text{var}[\Phi_{i,t}(\theta_j)] &= \frac{1}{t^2} \sum_{\tau=0}^{t-1} \sum_{k=1}^n \left[ \prod_{\rho=1}^{t-1-\tau} W(t-\rho) \right]_{ik}^2 \text{var}[\log \ell_k(s_k|\theta_j)] \\ &\leq \frac{1}{t} \sum_{k=1}^n \text{var}[\log \ell_k(s_k|\theta_j)].\end{aligned}$$

Hence, bounded variance assumption (A2) guarantees  $\Phi_{i,t}(\theta) - \mathbb{E}[\Phi_{i,t}(\theta)] \xrightarrow{p} 0$ . It can be shown [22] that

$$\mathbb{E}[W(t)] = \mathbb{E}[W(0)] = I - \frac{1}{2n}D + \frac{P + P^T}{2n},$$

where the diagonal matrix  $D$  is of the form  $D_i = \sum_{j=1}^n [P_{ij} + P_{ji}]$ . The fact that the sequence  $\{W(t)\}_{t=0}^\infty$  is i.i.d. and doubly stochastic, and the second largest eigenvalue of  $\mathbb{E}[W(t)]$  is less than one in magnitude (A3), entails [25]

$$W(t)W(t-1)\dots W(1) \longrightarrow \frac{1}{n}\mathbf{1}\mathbf{1}^T,$$

almost surely, which results in

$$\begin{aligned}\mathbb{E}[\Phi_{i,t}(\theta)] &= \sum_{k=1}^n \frac{1}{t} \sum_{\tau=0}^{t-1} \left[ \prod_{\rho=1}^{t-1-\tau} W(t-\rho) \right]_{ik} \mathbb{E}_{s_k}^*[\log \ell_k(s_k|\theta)] \\ &\longrightarrow \frac{1}{n} \sum_{k=1}^n \mathbb{E}_{s_k}^*[\log \ell_k(s_k|\theta)],\end{aligned}$$

where we used the fact that Cesàro mean preserves the limit. ■

There is an interesting connection between the previous lemma and the distributed MAP algorithm proposed in [7] where authors establish that the point maximizer of  $\Phi_{i,t}(\theta)$  over  $\Theta$  converges in probability to  $\theta^*$  in a strongly connected network. However, we still need to demonstrate that the estimator  $\mu_{i,t}(\theta)$  is weakly consistent for any agent  $i$ . To this end, we prove that applying a Bayes-like update in the same time-scale of receiving  $z_{i,t}$ , the belief vector  $\mu_{i,t}(\theta)$  converges in probability to the unique maximizer of (1) which is a Dirac distribution over  $\theta^*$ .

**Theorem 5:** Given conditions in the Lemmas 3 and 4, agent  $i$ 's estimator is weakly consistent, that is,

$$\mu_{i,t}(\theta) \xrightarrow{p} \mu^* = \mathbf{e}_1 \text{ as } t \rightarrow \infty.$$

*Proof:* We have the explicit form of the estimator  $\mu_{i,t}$  according to Lemma 3. Therefore,

$$\begin{aligned}\mu_{i,t}(\theta^* = \theta_1) &= \frac{\mu_{i,0}(\theta_1) \exp[t\Phi_{i,t}(\theta_1)]}{\sum_{j=1}^m \mu_{i,0}(\theta_j) \exp[t\Phi_{i,t}(\theta_j)]} \\ &= \left( 1 + \sum_{j \geq 2} \frac{\mu_{i,0}(\theta_j)}{\mu_{i,0}(\theta_1)} \exp[t\Phi_{i,t}(\theta_j) - t\Phi_{i,t}(\theta_1)] \right)^{-1}.\end{aligned}$$

Under purview of Lemma 4 and equation (3),  $[\Phi_{i,t}(\theta_j) - \Phi_{i,t}(\theta_1)]$  converges to a negative number for any  $j \geq 2$ , and hence  $\mu_{i,t}(\theta_1) \xrightarrow{p} 1$ . The fact that  $\mu_{i,t}(\theta) \in \Delta\Theta$  implies that  $\mu_{i,t}(\theta_j) \xrightarrow{p} 0$  for all  $j \geq 2$ , so  $\mu_{i,t}(\theta) \xrightarrow{p} \mu^* = \mathbf{e}_1$ . ■

Theorem 5 too underscores the trade-off between the adaptation and learning in the network. In many distributed optimization settings the stepsize sequence must vanish to

allow nodes to reach consensus. However, the result of Theorem 5 holds for unit stepsize sequence which guarantees learning even under continuous information injection to the network. This stems from the fact that the algorithm allows  $z_{i,t}$  to grow unboundedly in each direction, while it lets the true state to be the influential component by having the largest exponential rate in the generated Gibbs distribution.

### B. Learning Rate Analysis

In this section we characterize the convergence rate of the estimator  $\mu_{i,t}(\theta)$ . More specifically, we prove that convergence occurs exponentially fast with a rate dependent on the *average expected discrimination information* for  $\theta_1 = \theta^*$  over  $\theta_2$ , where  $\theta_2$  is the state with the second largest expected log-likelihood (3).

**Definition 6:** The expected discrimination information of agent  $i$  for  $\theta_1 = \theta^*$  over any  $\theta_j$  is

$$D_{KL}(\ell_i(\cdot|\theta_1) \parallel \ell_i(\cdot|\theta_j)) = \mathbb{E}_{s_i}^* \left[ \log \frac{\ell_i(s_i|\theta_1)}{\ell_i(s_i|\theta_j)} \right].$$

Denoting by  $D(\theta_j)$ , the *average expected discrimination information* for  $\theta_1 = \theta^*$  over  $\theta_j$  is defined as

$$D(\theta_j) \triangleq \frac{1}{n} \sum_{i=1}^n D_{KL}(\ell_i(\cdot|\theta_1) \parallel \ell_i(\cdot|\theta_j)). \quad (20)$$

As an immediate consequence of the definition above, one can see from Lemma 4 that  $D(\theta_j) = \Phi_\infty(\theta_1) - \Phi_\infty(\theta_j)$  for any  $j \geq 1$ , and  $D(\theta_1) = 0$ .

**Theorem 7:** Given conditions in the Lemmas 3 and 4, for any  $\epsilon > 0$  and  $t$  large enough, the estimator  $\mu_{i,t}(\theta_1)$  can be bounded as

$$|\mu_{i,t}(\theta_1) - 1| \leq \mathcal{K} \exp[(-D(\theta_2) + \epsilon)t], \quad (21)$$

with probability at least  $1 - \delta(\epsilon, t)$ , where  $\mathcal{K}$  is a constant.

*Proof:* Following the lines in the proof of Theorem 5, we have

$$\begin{aligned}\mu_{i,t}(\theta_1) &= \left( 1 + \sum_{j \geq 2} \frac{\mu_{i,0}(\theta_j)}{\mu_{i,0}(\theta_1)} \exp[t\Phi_{i,t}(\theta_j) - t\Phi_{i,t}(\theta_1)] \right)^{-1} \\ &\geq 1 - \sum_{j \geq 2} \frac{\mu_{i,0}(\theta_j)}{\mu_{i,0}(\theta_1)} \exp[t\Phi_{i,t}(\theta_j) - t\Phi_{i,t}(\theta_1)],\end{aligned}$$

where in the last step we used the inequality

$$1 - \lambda \leq (1 + \lambda)^{-1} \quad \forall \lambda \geq 0.$$

Letting  $b_j \triangleq \Phi_{i,t}(\theta_j) - \Phi_\infty(\theta_j)$ , we derive

$$\begin{aligned}|\mu_{i,t}(\theta_1) - 1| &\leq \sum_{j \geq 2} \frac{\mu_{i,0}(\theta_j)}{\mu_{i,0}(\theta_1)} \exp[t\Phi_{i,t}(\theta_j) - t\Phi_{i,t}(\theta_1)] \\ &\leq \max_k \frac{\mu_{i,0}(\theta_k)}{\mu_{i,0}(\theta_1)} \sum_{j \geq 2} \exp[t\Phi_{i,t}(\theta_j) - t\Phi_{i,t}(\theta_1)] \\ &= \max_k \frac{\mu_{i,0}(\theta_k)}{\mu_{i,0}(\theta_1)} \sum_{j \geq 2} \exp[(-D(\theta_j) + b_j - b_1)t].\end{aligned}$$

One can see in the proof of Lemma 4 that  $\text{var}[\Phi_{i,t}(\theta_j)]$  decays with a rate  $C/t$  for some constant  $C > 0$ ; hence, for any  $\epsilon > 0$  and  $j \geq 1$ , by Chebyshev's inequality we obtain

$$\mathbb{P}(|b_j| \geq \epsilon) \leq \frac{C}{\epsilon^2 t}.$$

Combining with  $D(\theta_m) \geq \dots \geq D(\theta_2) > D(\theta_1) = 0$  by (3), for any  $\epsilon > 0$  and  $t$  large enough, we have

$$\begin{aligned} |\mu_{i,t}(\theta_1) - 1| &\leq \max_k \frac{\mu_{i,0}(\theta_k)}{\mu_{i,0}(\theta_1)} \sum_{j \geq 2} \exp[(-D(\theta_2) + 2\epsilon)t] \\ &= (m-1) \max_k \frac{\mu_{i,0}(\theta_k)}{\mu_{i,0}(\theta_1)} \exp[(-D(\theta_2) + 2\epsilon)t], \end{aligned}$$

with probability at least  $1 - \frac{C}{\epsilon^2 t}$ . Hence, the constants in (21) are determined as

$$\mathcal{K} = (m-1) \max_k \frac{\mu_{i,0}(\theta_k)}{\mu_{i,0}(\theta_1)} \quad \text{and} \quad \delta(\epsilon, t) = \frac{4C}{\epsilon^2 t},$$

and we are done.  $\blacksquare$

Theorem 7 suggests that the proposed distributed stochastic learning method in (13)-(14)-(15) converges exponentially fast with high probability. Moreover, agents learn the true state with a rate dependent on the KL divergence of observations under the true state from observations under the second likeliest state. This, indeed, stresses the efficiency of the algorithm.

## V. CONCLUSION

We studied a distributed parameter estimation problem over networks when agents receive a sequence of i.i.d. signals but the signals are not informative enough to identify the true parameter. Using a randomized, gossip dual averaging, agents aggregate local log-likelihood functions, and then perform a Bayes-like update on the averaged information to collectively recover the truth. Assuming connectivity of the network and global identifiability of the true state, we showed that agents beliefs reach consensus and collapse to a degenerate distribution over the true parameter, and with high probability the convergence is exponentially fast. We also proved that the rate of exponential depends on the KL divergence of observations under true state from observations under second likeliest state. As a salient feature of the algorithm, we showed that contrary to other stochastic gradient descent methods, the stepsize can be chosen to be fixed and set to 1. Future directions include addition of dynamics to the parameter and relaxing the independence conditions on observations as well as specialization to the Gaussian case, where one only needs to update the mean and variance.

## ACKNOWLEDGMENTS

The authors would like to thank Robin Pemantle for many helpful comments and discussions.

## REFERENCES

- [1] V. Borkar and P. Varaiya, "Asymptotic agreement in distributed estimation," *IEEE Transactions on Automatic Control*, vol. 27, no. 3, pp. 650–655, 1982.
- [2] J. N. Tsitsiklis, "Decentralized detection by a large number of sensors," *Mathematics of Control, Signals, and Systems*, vol. 1, no. 2, pp. 167–182, 1988.
- [3] J. N. Tsitsiklis *et al.*, "Decentralized detection," *Advances in Statistical Signal Processing*, vol. 2, pp. 297–344, 1993.
- [4] E. Mossel and O. Tamuz, "Efficient bayesian learning in social networks with gaussian estimators," *Arxiv preprint arXiv:1002.0747*, 2010.
- [5] U.A.Khan, S. Kar, A. Jadbabaie, and J. Moura, "On connectivity, observability, and stability in distributed estimation," in *49th IEEE Conference on Decision and Control (CDC)*. IEEE, 2010, pp. 6639–6644.
- [6] S. Kar, J. Moura, and K. Ramanan, "Distributed parameter estimation in sensor networks: Nonlinear observation models and imperfect communication," *IEEE Transactions on Information Theory*, vol. 58, no. 6, pp. 3575–3605, 2012.
- [7] K. Rad and A. Tahbaz-Salehi, "Distributed parameter estimation in networks," in *49th IEEE Conference on Decision and Control (CDC)*. IEEE, 2010, pp. 5050–5055.
- [8] A. Jadbabaie, P. Molavi, A. Sandroni, and A. Tahbaz-Salehi, "Non-bayesian social learning," *Games and Economic Behavior*, vol. 76, no. 1, pp. 210–225, 2012.
- [9] O. Dekel, R. Gilad-Bachrach, O. Shamir, and L. Xiao, "Optimal distributed online prediction using mini-batches," *The Journal of Machine Learning Research*, vol. 13, pp. 165–202, 2012.
- [10] J. Tsitsiklis, "Problems in decentralized decision making and computation." DTIC Document, Tech. Rep., 1984.
- [11] A. Jadbabaie, J. Lin, 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.
- [12] M. Mesbahi and M. M. Egerstedt, *Graph theoretic methods in multi-agent networks*. Princeton Univ Press, 2010.
- [13] F. Bullo, J. Cortés, and S. Martínez, *Distributed control of robotic networks: a mathematical approach to motion coordination algorithms*. Princeton Univ Pr, 2009.
- [14] R. Olfati-Saber and J. Shamma, "Consensus filters for sensor networks and distributed sensor fusion," in *44th IEEE Conference on Decision and Control*, Seville, Spain, Dec. 2005, pp. 6698 – 6703.
- [15] A. Nedic and A. Ozdaglar, "Distributed subgradient methods for multi-agent optimization," *IEEE Transactions on Automatic Control*, vol. 54, no. 1, pp. 48–61, 2009.
- [16] I. Lobel and A. Ozdaglar, "Distributed subgradient methods over random networks," in *Proc. Allerton Conf. Commun., Control, Comput*, 2008.
- [17] S. Ram, A. Nedic, and V. Veeravalli, "Distributed stochastic subgradient projection algorithms for convex optimization," *Journal of optimization theory and applications*, vol. 147, no. 3, pp. 516–545, 2010.
- [18] A. Nedic, A. Olshevsky, A. Ozdaglar, and J. Tsitsiklis, "On distributed averaging algorithms and quantization effects," *IEEE Transactions on Automatic Control*, vol. 54, no. 11, pp. 2506–2517, 2009.
- [19] C. Lopes and A. Sayed, "Incremental adaptive strategies over distributed networks," *IEEE Transactions on Signal Processing*, vol. 55, no. 8, pp. 4064–4077, 2007.
- [20] J. Duchi, A. Agarwal, and M. Wainwright, "Dual averaging for distributed optimization: convergence analysis and network scaling," *IEEE Transactions on Automatic Control*, pp. 592–607, March 2012.
- [21] Y. Nesterov, "Primal-dual subgradient methods for convex problems," *Mathematical programming*, vol. 120, no. 1, 2009.
- [22] S. Boyd, A. Ghosh, B. Prabhakar, and D. Shah, "Randomized gossip algorithms," *IEEE Transactions on Information Theory*, vol. 52, no. 6, pp. 2508–2530, 2006.
- [23] T. M. Cover and J. A. Thomas, *Elements of information theory*. John Wiley & Sons, 2012.
- [24] D. Blackwell and L. Dubins, "Merging of opinions with increasing information," *The Annals of Mathematical Statistics*, vol. 33, no. 3, pp. 882–886, 1962.
- [25] A. Tahbaz-Salehi and A. Jadbabaie, "Consensus over ergodic stationary graph processes," *IEEE Transactions on Automatic Control*, vol. 55, no. 1, pp. 225–230, 2010.